

PCT

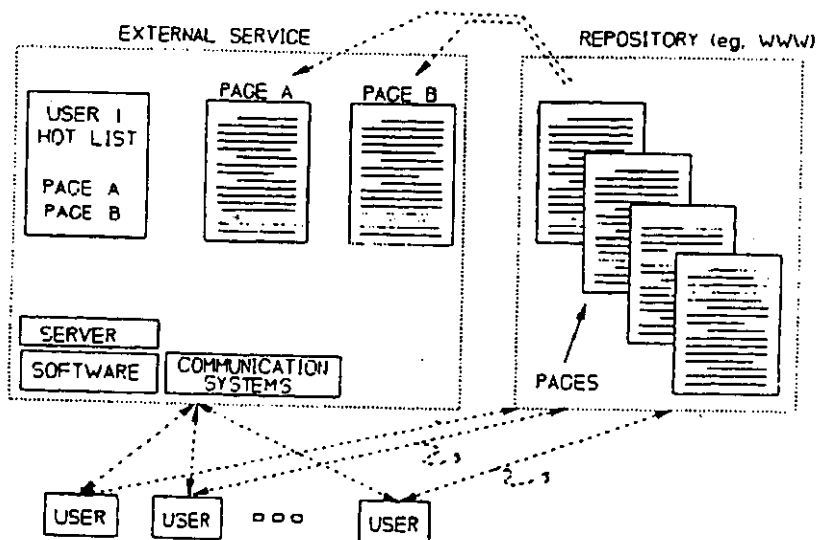
WORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau



INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification 6: G06F 17/30	AI	(11) International Publication Number: WO 97/15890
		(43) International Publication Date: 1 May 1997 (01.05.97)
(21) International Application Number: PCT/US96/17142		(81) Designated States: CA, JP, European patent (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE). Published With international search report.
(22) International Filing Date: 25 October 1996 (25.10.96)		
(30) Priority Data: 08/549,359 27 October 1995 (27.10.95) US		
(71) Applicant: AT & T CORP. [US/US]; 32 Avenue of the Americas, New York, NY 10013-2412 (US).		
(72) Inventors: BALL, Thomas, J.; 1236 Whittingham Circle, Naperville, IL 60540 (US). DOUGLIS, Frederick; 248 Harlech Way, Somerset, NJ 08873 (US).		
(74) Agent: DWORETSKY, Samuel, H.; c/o H.T. Brendzel, AT & T Corp., P.O. Box 4110, Middletown, NJ 07748 (US).		

(54) Title: IDENTIFYING CHANGES IN ON-LINE DATA REPOSITORIES



(57) Abstract

A system for accessing documents contained in a remote repository, which change in content from version-to-version. The system allows users to specify lists of documents of interest. Based on the lists, the system maintains an archive, which contains a copy of one version of each listed document, and material from which the other versions can be reconstructed. The system periodically compares the archive with current versions of the documents located in the repository, and updates the archive, thereby maintaining the ability to reconstruct current versions. The system also monitors access to the versions by each user. When a user calls for a current version, the system presents the current version, and indicates what parts of the current version have not been previously accessed by the user.

BEST AVAILABLE COPY

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AM	Armenia	GB	United Kingdom	MW	Malawi
AT	Austria	GE	Georgia	MX	Mexico
AU	Australia	GN	Guinea	NE	Niger
BB	Barbados	GR	Greece	NL	Netherlands
BE	Belgium	HU	Hungary	NO	Norway
BF	Burkina Faso	IE	Ireland	NZ	New Zealand
BG	Bulgaria	IT	Italy	PL	Poland
BJ	Benin	JP	Japan	PT	Portugal
BR	Brazil	KE	Kenya	RO	Romania
BY	Belarus	KG	Kyrgyzstan	RU	Russian Federation
CA	Canada	KP	Democratic People's Republic of Korea	SD	Sudan
CF	Central African Republic	KR	Republic of Korea	SE	Sweden
CG	Congo	KZ	Kazakhstan	SG	Singapore
CH	Switzerland	LI	Liechtenstein	SI	Slovenia
CI	Côte d'Ivoire	LK	Sri Lanka	SK	Slovakia
CM	Cameroon	LR	Liberia	SN	Senegal
CN	China	LT	Lithuania	SZ	Swaziland
CS	Czechoslovakia	LU	Luxembourg	TD	Chad
CZ	Czech Republic	LV	Latvia	TG	Togo
DE	Germany	MC	Monaco	TJ	Tajikistan
DK	Denmark	MD	Republic of Moldova	TT	Trinidad and Tobago
EE	Estonia	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	UG	Uganda
FI	Finland	MN	Mongolia	US	United States of America
FR	France	MR	Mauritania	UZ	Uzbekistan
GA	Gabon			VN	Viet Nam

IDENTIFYING CHANGES IN ON-LINE DATA REPOSITORIES

The invention concerns presentation of a current version of a document retrieved from a data repository. The presentation indicates changes made in the document since the viewer accessed a previous version.

BACKGROUND OF THE INVENTION

Information which is stored in computerized systems can change frequently, and without notice. As an example, software under development frequently involves many persons, and is commonly stored at a central location. Each person can change the software on an ad hoc basis, without knowledge of others.

In such systems containing changeable data, a person who examines information on a given day does not, in general, know whether, and how, the information has changed since a previous examination. Consequently, the person must spend time comparing currently available information with previous versions of the information.

Software exists for facilitating this comparison. For example, systems known as "version control systems," or "revision control systems," store data which represents multiple versions of different documents, as indicated in Figure 1A. In that Figure, the DATA is indicated, together with dashed loops which indicate the VERSIONS.

The loops indicate that the VERSIONS are contained in, and

WO 97/15890

PCT/US96/17142

MISSING UPON TIME OF PUBLICATION

WO 97/15890

PCT/US96/17142

MISSING UPON TIME OF PUBLICATION

interspersed among a much larger document. Small arrows point to changes, which are primarily additions in this case. The change in the "last update" date give an example of text being replaced. Here the page's author had highlighted the changes manually with small icons as well. The banner at the top of the page was inserted by HTMLDIFF.

Figure 12 illustrates version histories which give the user a chance to compare any two versions, or to go directly to a selected version.

Figure 13 illustrates output of W3NEWER, and shows a number of anchors (the descriptive text originates from the hot list). The anchors marked "changed" have modification dates after the time which the user's browser history indicates the URL was last seen. Some URLs were not checked at all, and others were checked and are known to have been seen by the user.

Figure 14 demonstrates use of a SNAPSHOT facility, which allows a user to specify an operation on a URL. In this example, DOUGLIS@RESEARCH.ATT.COM is "remembering" URL HTTP://SNAPPLE.CS.WASHINGTON.EDU:600/MOBILE/.

DETAILED DESCRIPTION OF THE INVENTION

A TECHNICAL APPENDIX, which is located at the end of this Specification, describes the invention in detail. Following the TECHNICAL APPENDIX is a COMPUTER PROGRAM LISTING, which contains code which implements one form of the invention.

An illustrative embodiment of the invention is given in the

discussion below.

Overview of Invention

A commonly used repository of information is known as the World Wide Web, or WWW. In the WWW, providers of information make their information available to users in the form of "pages." Each page is assigned a name, which distinguishes the page from other pages, and allows a user to locate the page.

The WWW provides information using an information retrieval-and-display approach called "hypertext." In hypertext, a page may contain references to other pages, or other documents. A user can call up a page which is referenced, by clicking on the reference (called a URL, or Universal Resource Locator) with a pointing device. Figure 1B provides an example.

In Figure 1B, a document D is displayed to a user. References R refer to other documents. For example, R1 refers to D1, R2 refers to D2, and so on. The referenced documents themselves may contain their own references to other documents, such as R4, which refers to D4.

A user can retrieve a referenced document D, by clicking on the reference R which refers to it. For example, clicking on R1 causes retrieval and display of D1.

Under the invention, a user of the WWW initially identifies pages of interest. Document D in Figure 1B represents one page. These selected pages form a "hot list." Then, the invention does the following:

(a) Copies the hot-listed pages into an archive, which is a storage location separate from the WWW, and under independent control. After the copying, the original pages continue to reside in the WWW, and copies reside in the archive.

(b) Monitors, at later times, the original pages for changes, and archives the changes.

(c) Records the times when the user later accesses each hot-listed page.

(d) Whenever the user accesses a hot-listed page, presents the user with

- i) the current version of the page (which may differ from the initial copy which was stored in the archive); and

- ii) an option to compare selected versions of the page. The comparison is presented by performing a differencing operation on pairs of versions.

e) As an option, the invention also implements the steps described above with respect to documents referenced by the page. For example, in Figure 1A, if a user is viewing document D, the invention can present the current version of reference document D2, together with a history of D2.

More Detailed Description

Hot-List Pages are Stored in EXTERNAL SERVICE

Figure 1 illustrates a REPOSITORY of information, such as the WWW. For assistance in accessing the REPOSITORY, the invention provides the EXTERNAL SERVICE which includes:

- (a) SOFTWARE, such as that provided in the COMPUTER PROGRAM LISTING herein,
- (b) a SERVER, or other computer, which runs the software, and
- (c) COMMUNICATION SYSTEMS which link with both the users and the REPOSITORY.

The SERVER and the COMMUNICATION SYSTEMS located within the EXTERNAL SERVICE are known in the art. As indicated in the Figure, the EXTERNAL SERVICE is distinct from the REPOSITORY, and under separate control.

The invention does not disrupt the users' normal interaction with the REPOSITORY; the users can interact with both the REPOSITORY, as usual, and also with the EXTERNAL SERVICE. Dashed arrows 3 indicate the interaction. Several examples will provide illustrative modes of operation of the invention.

Example: Single User

Operation with respect to a single user will first be explained. Figure 2 shows a hot list 4, submitted by USER 1, which identifies pages A and B as being of interest to USER 1. The

invention allows the user to modify the hot list at later times. In response to the hot list, the invention copies pages A and B from the REPOSITORY, as indicated by the dashed arrows. These PAGES will be termed "base pages." At this time, the originals of PAGES A and B remain in the REPOSITORY, and copies reside in the EXTERNAL SERVICE.

Then, the invention periodically examines the originals of PAGES A and B, located in the REPOSITORY, for changes. In looking for changes, the invention first performs a preliminary check, based on information such as (1) dates of modification and (2) checksums.

Dates of modification may be added to a PAGE by the PAGE provider. These dates directly indicate whether the originally archived version has changed.

Checksums are generated by the invention. An example of a checksum is the numerical sum of all characters in a line, or on a page. If a checksum changes (indicating that the number of characters has changed), the change indicates a high probability that a change has occurred in the PAGE. (In practice, the checksums used are more complex than this simple example illustrates. Checksums are known in the art.)

If the preliminary check, either by dates of modification or checksums, indicates that changes have occurred, then the invention copies the present version of the PAGE into the EXTERNAL SERVICE, and compares it with the base page, in order to locate the changes. Computer programs for detecting such changes are known in the art,

and some examples are given in the TECHNICAL APPENDIX. A preferred program, not known in the prior art, is entitled W3NEWER, and was developed by the inventors. W3NEWER is contained in the listing located at the end of this Specification.

When changes are found, the invention stores them in the EXTERNAL SERVICE. Figure 3 illustrates storage of the changes, by the small boxes 6 located below PAGEs A and B. The DATEs within the boxes 6 indicate the dates on which the changes were saved.

Figure 3A illustrates how the invention displays the history of versions. Column 7 indicates the number assigned to each version by the invention. Column 8 indicates the times when the respective versions were retrieved by the invention. Column 8A allows a user to select a pair of versions for a differencing operation, as discussed below.

For ease of explanation, Figure 3 illustrates storage of base pages, which are early versions of PAGEs, together with subsequent changes, indicated by the boxes 6. However, in practice, it can be more efficient to perform storage in a reversed sense, by storing the latest version as the base page (instead of the early version) and storing the changes 6 from which early versions can be reconstructed. One reason is that users are expected to call for latest versions more frequently than early versions. Storage of the entire latest versions eliminates the need to reconstruct them.

The changes, together with their base pages, form an archive, which allows reconstruction of a PAGE as of any date desired. For

example:

-- PAGE A itself (ie, the base page), plus the changes labeled DATE 1, allow reconstruction of the version of PAGE A, as of DATE 1.

-- PAGE A itself, plus the changes labeled DATE 1 and DATE 2, allow reconstruction, as of DATE 2, and so on.

When USER 1 wishes to view PAGE A, the invention ordinarily retrieves and presents the current version. The invention also provides an option for reconstructing the PAGE, as of a date specified by the user, and presents it in the format shown in Figure 4. The program HTMLDIFF, contained in the listing, generates the image shown in Figure 4. The content of the page can be divided into three classes.

The first class contains material which has not changed. This class of material is displayed in the font, size, color, and background, as customary in documents downloaded from the REPOSITORY.

The second class represents changes, and contains material not present in the base page, but which has been added. Brackets 9 indicate such material. (The brackets 9 are part of Figure 4, and are not necessarily part of the page generated by the invention.) This material is presented in a particular font, particular size, particular color, and particular background. The choice of these parameters can be varied but, in general, they should be chosen to maximize contrast with the first class of material. In addition to the formatting described immediately

above, the added material is further highlighted by arrows 7.

The third class contains material which was deleted from the base page. Deleted material can be handled in at least three ways. One, deleted material can be simply deleted, so that the page presented to the reader contains no reference to the deleted material.

Two, the deleted material can be deleted, but a reference indicating the deletion is added, such as the phrase "Deleted material occurs here." In this case, the user can be given the option of fetching the deleted, non-visible, material.

Three, deleted material can be presented, but indicated as deleted, as by "redline" format, in which a horizontal line, perhaps red in color, is drawn through the deleted material.

Figure 3B illustrates a display, generated by the invention, which indicates which PAGES on a user's hot list have undergone changes.

Second Example: Multiple Users

In actual practice, multiple users are expected to use the invention. Each of them submits a hot list. In one approach of the invention, the procedure undertaken for a single user (described above) is repeated for multiple users: all PAGES, on all hot lists, are copied into the EXTERNAL SERVICE. Then, for each hot list, the originals of the PAGES, located within the REPOSITORY, are monitored for changes, and the changes are retrieved into the EXTERNAL SERVICE, as described above.

However, this approach contains inefficiencies. For example, a given PAGE will probably be identified by more than one hot list. Repeatedly copying that PAGE, for each hot list, would entail storage of multiple copies of the same PAGE. Further, repeatedly comparing the multiple copies with their originals in the REPOSITORY represents a waste of computer time: a single comparison would suffice. The invention reduces these inefficiencies by the approach shown in Figure 5.

This Figure represents a modification of Figure 4, to which a hot list for USER 2 has been added. The added hot list specifies PAGES A and C.

To process the new hot list, the invention first checks whether the PAGES identified on the added hot list are archived within the EXTERNAL SERVICE. Since PAGE A, plus its changes, are already contained within the archive, that PAGE is not copied. But PAGE C, which is not present in the ARCHIVE, is archived, as indicated by the dashed arrow.

At this time, all PAGES identified on all hot lists are contained within the archive. To emphasize this fact, PAGE A is indicated twice: once for USER 1, and a second time by a dashed page 14, for USER 2, although, as stated above, PAGE A is stored only once.

After archiving all necessary PAGES, the originals, located within the REPOSITORY, are periodically monitored for changes, as described above. The changes are copied to the archive of the EXTERNAL SERVICE.

Flow Chart

An exemplary flow chart is shown in Figure 6, which refers to a single-user case. In block 20, the EXTERNAL SERVICE accepts hot lists from users. Then, in block 23, the EXTERNAL SERVICE checks whether the PAGES identified on the hot lists are contained within the archive. If not, the PAGES are copied from the REPOSITORY, as indicated by block 26.

Then the logic proceeds to block 29, where the originals of the PAGES, located in the REPOSITORY, are examined for changes. The examination can include the preliminary checks (for checksums and dates of modification) discussed above. When changes are found, the entire PAGE containing them is downloaded to the EXTERNAL SERVICE, and the changes, indicated by blocks 6 in Figure 3, are derived. Block 32 indicates relevant information stored in the EXTERNAL SERVICE.

As users access the PAGES, block 35 monitors the times of the accesses, in order to identify which versions of each PAGE the user viewed last. These times are stored, as indicated by block 32 and dashed arrow 37. These times are used to determine which changes in Figure 4 are to be identified as new material, when a PAGE is called by each user. An example will illustrate.

Figure 7, top, illustrates the time-history of changes made to PAGE A. USER 1 accessed this PAGE at time 2, as indicated. Block 35 in Figure 6 monitors and records this time (at TIME 2 in Figure 7, and not earlier, of course).

If USER 1 again accesses the PAGE at time 5, then the invention presents VERSION 1 to the USER. However, if the user accesses the PAGE at time 11, VERSION 2 had been created since the last access by USER 1. The invention had previously identified the changes, and copied them as indicated in Figure 3. Now, at the access at time 11, the invention presents VERSION 1, plus the changes which make VERSION 2, because block 35 in Figure 6 indicates that the USER has not seen VERSION 2.

Returning to the flow chart of Figure 6, block 39 indicates that, when a USER calls for a PAGE, the invention presents the current version, and indicates the changes made (as in Figure 4) since the USER last accessed that page. In the example immediately above, the invention presents VERSION 2 of PAGE A, as in Figure 7, and indicates the changes made since VERSION 1, because VERSION 1 was the last accessed by USER 1.

The flow chart of Figure 6 should not be read as limiting the invention to a linear, sequential mode of operation. In practice, multiple users can present hot lists simultaneously, and other operations shown in the flow chart can also occur together.

Third Example: Notification of Changes

The invention can notify USERS when changes in their hot-listed PAGES occur, as indicated by the dashed block 40 in Figure 6. This notification can take the form of a flag which is associated with the BASE PAGE in Figure 8. When the USER logs into the EXTERNAL SERVICE, the invention notifies the USER of the

changes to the respective PAGES. Figure 3B illustrates one approach to identifying PAGES which have changed.

Other types of notification are possible. For example, the invention need not wait for a user to access a PAGE. The invention can notify the user when changes have been found, as by sending an electronic mail message to the user.

Fourth Example: Common Hot List

The invention can maintain a predetermined hot-list, for a community of users. This hot list contains a list of PAGES which are considered to be of general interest to the community. This hot list, and the PAGES identified on it, are made publicly available, to all users, but on a read-only basis. Users cannot modify the hot list, or the pages.

This predetermined hot list can serve as an instructional tool, to educate users in the operation of the invention, and to demonstrate desirable features.

One Architecture of Data Storage

An illustrative approach to storage of the information identified in block 32 of the flow chart of Figure 6 is illustrated in Figure 8, which is explained with reference to Figure 7.

Figure 7 illustrates hypothetical changes to the three PAGES identified by the two hot lists of Figure 5. PAGE A underwent changes at times 7 and 13. Page B underwent changes at time 10, and so on.

In Figure 8, the arrows extending from the symbols "USER 1", etc., indicate the times of access by the users. For example, USER 1 accessed PAGE A, VERSION 1, at time 2. USER 1 then accessed PAGE A, VERSION 2, at time 9, and so on.

The invention maintains a TABLE of these times, as indicated on the right side of Figure 8, together with a list of PAGES, or documents, owned by each USER. Ownership is determined by the hot lists. The invention also maintains (a) the BASE PAGES, (b) the changes to each, and (c) the times of each change, as indicated on the left side of the Figure. From this data, the invention is able to reconstruct any PAGE, as of any date subsequent to the date of the BASE PAGE.

Additional Considerations

1. One definition of "page" is that it refers to a unit of data, stored in a system, which is identified by a specific name. (In the WWW, all pages have unique names.) Other terms can refer to such units of data, such as "files" and "documents." In general, the particular name used will depend on the system storing the data.

2. One definition of "repository" is a collection of data, which is accessible by computer. The repository may be available to the public, or access may be limited. In general, repositories are expected to be distributed, meaning that the storage locations are physically distributed over a wide geographic area, and linked

together by a communication system.

3. It was stated above that the invention can reconstruct a page as of any selected date. The reconstruction is based on the changes 6 in Figure 3. These changes are detected periodically, and the periodicity is determined by each user of the system, subject to limits imposed by the designer and system administrator.

For example, user A can specify a period of one day for checking for changes in the pages on user A's hot list; user B can specify a different period for B's pages, such as one week. The system administrator can specify that no period, for any user, can be shorter than one hour.

Consequently, changes in a page, located in the REPOSITORY, will only appear in a reconstruction done by the EXTERNAL SERVICE after the changes have been detected, and not earlier. An example will illustrate this distinction.

Assume that the invention looks for changes on odd-numbered dates. Thus, a change occurring on the fourth of a month will be detected on the fifth. However, if a user happens to call for reconstruction on the fourth, the change occurring on the fourth will not appear in the reconstruction. Only changes occurring as of the prior detection, namely, as of the third, will appear.

It is expected that the detection process will be performed sufficiently often that the influence of this factor will be negligible.

4. The invention can extend its differencing function (ie, the examination of pages for changes) to pages referenced by the page accessed by the user. For example, if the user accesses document D in Figure 1B, the invention can detect changes in all documents referenced by document D, such as D1, D2, and D3.

In another embodiment, the differencing can extend to the documents which are, in turn, referenced by the referenced documents. For example, the referenced documents (D1, D2, and D3) refer to D5 and D6. These latter documents (D5 and D6) can be differenced also, as can be the documents which they reference, and so on.

5. The invention provides all information from which a current version of a PAGE may be derived. Figure 4 gives an example. Figure 4 contains all such information, together with other information which indicates changes since a previous version.

6. The discussion above presumed that comparison, or differencing, between different versions of a PAGE was done within the EXTERNAL SERVICE. This is not strictly necessary; the comparison can be done at any convenient location. Further, the preliminary checking for the existence of changes can be done at any convenient location.

7. In data storage systems, names are given to the units of information (e.g., documents, pages, records), although the names

can be different in different databases. However, the names of the units, in general, remain the same throughout time, despite changes which are made to the information contained in the unit. Therefore, one definition of the term "version" refers to a unit of information, which is different from a previous unit of the same name.

8. The REPOSITORY in Figure 1 is, in general, located remotely from the EXTERNAL SERVICE. Communication is undertaken by any convenient approach, such as a public-access communication network known as the INTERNET.

In general, the REPOSITORY is under independent control of the EXTERNAL SERVICE. One ramification of this independent control is that the type of processing done to the PAGES copied into the EXTERNAL SERVICE is controlled by the EXTERNAL SERVICE, and not by the REPOSITORY. For example, (a) the particular processes used in locating and storing differences, (b) the frequency of processing, and (c) the mode of notifying a user, are controlled by the designer of the EXTERNAL SERVICE. The operator of the REPOSITORY has no involvement in this processing.

9. Figure 9 illustrates another form of the invention. The invention maintains base pages 30 within the EXTERNAL SERVICE, as required by the hot lists 36. The base pages 30 were downloaded from respective repositories 42A, 42B, etc.

The invention periodically monitors the originals 30A of the

pages, located in the repository 42, for changes, and stores the changes within the EXTERNAL SERVICE. The invention notifies users when changes are found in pages on their hot lists (notification is not shown).

A version control system 39 allows users to fetch and view any version of any page.

10. The different versions of documents may contain drawings, files from which sound may be generated, files which produce video clips and animation, and other components which do not consist strictly of alphanumeric characters. The invention detects the existence of changes in such components, and marks the existence of the changes, in the display as shown in Figure 4, without necessarily identifying in detail the nature of the changes.

11. A primary use of the invention is envisioned in the situation shown in Figure 10. The EXTERNAL SERVICE obtains copies of PAGES from a REPOSITORY, such as WWW. However, the EXTERNAL SERVICE is given no authority to replace or modify the pages contained in the REPOSITORY. To the EXTERNAL SERVICE, the PAGES represent read-only data, as indicated by the "X" over arrow 50, which indicates a write operation.

The EXTERNAL SERVICE performs differencing between currently copied versions of pages, and DATA representing previous versions. The DATA stored in the EXTERNAL SERVICE can be both read, and written to, by the EXTERNAL SERVICE. The EXTERNAL SERVICE

reconstructs any version on demand, and also indicates differences between any two versions selected by a user, as discussed above. These functions can be accomplished by a prior-art Revision Control System, RCS (also called a Version Control System), or by the code contained in the listing contained in this Specification.

12. In one form of the invention, the PAGEs retrieved are written in a "markup language," such as HyperText Mark-up Language (HTML). A mark-up language, in general, contains two types of codes, interspersed among the actual text of a document.

One type indicates how the PAGEs are to be displayed. For example, some codes indicate paragraph indentation, other codes indicate font styles, yet other codes indicate style of font, within a font, such as italicizing, underlining, double-striking, or bold printing. This type of code is referred to as format-defining.

A second type of code can identify an image, such as a bit-mapped file located elsewhere. When such a code is read by the system displaying the PAGE, a copy of the image is retrieved, and displayed within the PAGE, at the location specified by the code. This type of code is referred to as content-defining.

The invention does not treat changes in the format-defining codes as changes in content. Thus, a PAGE which changes in layout, or typestyle, only, is not designated as a changed page.

The differencing program contained in the COMPUTER PROGRAM LISTING compares different versions on a subunit-by-subunit basis.

For example, the program compares corresponding sentences in different versions, and the sentences are detected by sentence terminators. (Longer subunits can be used, such as paragraphs or pages.) The sentence terminators are a subset of the markup language. Specifically, the terminators are format-defining codes.

COMPUTER PROGRAM LISTING

The program listing is divided into three sections.

1. HTMLDIFF, comprising:

- html_diff.sml (5 pages),
- diff.sml (3 pages),
- mlweb.sml (4 pages), and
- html.lex (one page).

2. W3NEWER (17 pages).

3. NOHANDS, comprising:

- nohandsBE (11 pages),
- no-hands.cgi (3 pages),
- rcsdiff.cgo (4 pages), and
- snapshot.cgi (3 pages).

NOHANDS is an overall program set which utilizes W3NEWER and HTMLDIFF.

TECHNICAL APPENDIX

A TECHNICAL APPENDIX, totalling 12 pages, two of which are

blank, follows, and refers to Figures 11 - 14.

Numerous substitutions and modifications can be undertaken without departing from the true spirit and scope of the invention. What is desired to be secured by Letters Patent is the invention as defined in the following claims.

CLAIMS

1. A system, comprising:
 - a) means for copying versions of pages from an external repository, which provides pages to the system on a read-only basis;
 - b) means for storing data from which
 - i) selected versions can be reconstructed;
 - and
 - ii) differences between selected versions can be identified.
2. A system according to claim 1, in which the versions reconstructed are specified on lists provided by users.
3. In a system for allowing multiple users to interface with another system which contains versions of data, at least some of which change as time progresses, the improvement comprising:
 - a) means for allowing a user to identify a current version and a previous version; and
 - b) means for displaying the current version in a format which distinguishes the current version from the previous version.
4. A method of communicating with a repository which stores

data units in the form of current versions which change over time, comprising:

- a) accepting lists of data units from multiple users;
- b) maintaining an archive, which contains material which allows reconstruction of selected versions of data units contained on the lists; and
- c) at intervals, comparing current versions with archived material, and updating the archive.

5. Method according to claim 4, and further comprising the step of

- d) recording which versions each user has accessed.

6. Method according to claim 5, and further comprising the steps of:

- e) accepting a request from a user for a data unit;
- f) presenting a current version of the data unit; and
- g) highlighting material within the current version which was absent from previous versions accessed by the user.

7. A method of communicating with a repository which stores data units in the form of versions which change over time, comprising:

- a) maintaining an archive, which includes material which allows reconstruction of selected versions of data units;

- b) accepting lists of data units from users;
- c) for each list accepted, checking the archive and, if listed data units are not archived, copying a current version from the repository;
- d) at intervals,
 - i) finding differences between (A) versions contained in the archive, and (B) current versions, contained in the repository, and
 - ii) storing the differences in the archive;
- e) maintaining information about each user, which indicates most recent versions accessed by the user; and
- f) accepting a request from a user for a current version and, in response,
 - i) copying the current version from the repository;
 - ii) finding differences between the current version and the most recent versions accessed; and
 - iii) displaying the current version in a format which highlights the differences.

8. Apparatus for comparing a first text written in a markup language, which includes sequences of words and markup commands, with a second such text comprising:

- a) means for comparing the first and second text, which compares sequences in the first and second text which

are defined by a group of terminators, including sentence-ending punctuation and markup commands, belonging to a first subset thereof, and produces a result which indicates differences between the texts; and
b) means for receiving the result and displaying the differences in response thereto.

9. Apparatus for comparing versions of pages in a repository, comprising:

- a) means for detecting that a second page is a new version of a first page;
- b) means responsive to the detecting means for comparing the first page with the second page to produce a result which indicates differences between the first and second pages; and
- c) means responsive to the result for displaying the differences.

10. A method for detecting whether a page in a repository has been modified since it was last examined by a given user, the method comprising the steps performed in the computer system used by the given user of:

- a) maintaining a first record of last viewed times for pages in that computer system which indicates, for each page in a first set of pages, the last time at which the given user viewed the page and for a given page for which

a last viewed time is recorded in the first record, performing the steps of obtaining a last modified time for the given page from a second record in that computer system of last modified times of pages and determining whether the last modified time from the second record is a time such that the time is later than the last viewed time for the given record;

b) if the second record does not provide a last modified time which is such a time, obtaining a last modified time from a source external to that computer system to which the user has access and determining whether the last modified time from the external source is such a time; and

c) if the last modified time obtained from the external source is such a time, updating the second record with that last modified time.

11. A method for maintaining a user-defined version history of a page, the method comprising the steps of:

a) receiving an indication from the user specifying that a given version of a page is to be retained in a version history for that page;

b) if the given version is the first version of the page in the version history, storing a copy of the given version in the version history; and

c) otherwise storing at least any difference between the

given version and the most recent previous version in
the version history.

TECHNICAL APPENDIX

Tracking and Viewing Changes on the Web

Abstract

We describe a set of tools that detect when World-Wide-Web pages have been modified and present the modifications visually to the user through marked-up HTML. The tools consist of three components: *w3newer*, which detects changes to pages; *snapshot*, which permits a user to store a copy of an arbitrary Web page and to compare any subsequent version of a page with the saved version; and *htmldiff*, which marks up HTML text to indicate how it has changed from a previous version. We refer to the tools collectively as the *Network-Oriented HTML Archival, Notification, and Differencing System* (NO HANDS). This paper discusses several aspects of NO HANDS, with an emphasis on systems issues such as scalability, security, and error conditions.

1 Introduction

Use of the World-Wide-Web (W^3) has increased dramatically over the past couple of years, both in the volume of traffic and the variety of users and content providers. The W^3 has become an information distribution medium for academic environments (its original motivation), commercial ones, and virtual communities of people who share interests in a wide variety of topics. Information that used to be sent out over electronic mail or USENET, both active media that go to users who have subscribed to mailing lists or newsgroups, can now be posted on a W^3 page. Users interested in that data then visit the page to get the new information.

The URLs of pages of interest to a user can be saved in a "holist" (known as a bookmark file in NetscapeTM), so they can be visited conveniently. How does a user find out when pages have changed? If users know that pages contain up-to-the-minute data (such as stock quotes), or are frequently changed by their owners, they may visit the pages often.

Other pages may be ignored, or browsed by the user only to find they have not changed.

In recent months, several tools have become available to address the problem of determining when a page has changed. One example of such a tool is *webwatch*, a product for WindowsTM that uses the HTTP HEAD command to find out when a page has been modified since it was last viewed by a user's web browser, and generates a report in HTML that allows the user to go directly to those updated pages. Another example is *w3new*, by Brooks Cutter, a public-domain perl script that runs on UNIX[®] [2].

Each of these tools suffers from a significant deficiency: while they provide the user with the knowledge that the page has changed, they do not show *how* the page has changed. Although a few pages are edited by their maintainers to highlight the most recent changes, often the modifications are not prominent, especially if the pages are large. Even pages with special highlighting of recent changes are problematic: if a user visits a page frequently, what is "new" to the maintainer may not be "new" to the user. Alternatively, a user who visits a page infrequently may miss changes that the maintainer deems to be old.

We have developed a system that efficiently tracks when pages change, compactly stores versions on a per-user basis, and automatically compares and presents the differences between pages. NO HANDS (*Network-Oriented HTML Archival, Notification, and Differencing System*) provides "personalized" views of versions of W^3 pages with three tools. The first, *w3newer*, is a more scalable version of Cutter's *w3new* modification tracking tool that periodically accesses the W^3 to find when pages on a user's hostlist have changed. The second, *snapshot*, allows a user to save versions of a page and later use a third tool, *htmldiff*, to see how it has changed. *Htmldiff* automatically compares two HTML pages and creates a "merged" page to show the differences with special HTML markups.

TM Netscape is a trademark of Netscape Communications Corporation.

TM Windows is a trademark of Microsoft.
[®] UNIX is a registered trademark of X/Open.

Tracking and Viewing Changes on the Web (NO HANDS)

checked. The threshold can vary depending on the URL, with *perl* pattern matching used to determine what threshold to apply. The first matching pattern is used. Table 1 gives an example of a *w3newer*.thresholds configuration file. Thresholds are specified as combinations of days (d) and hours (h), with 0 indicating that a page should be checked on every run of *w3newer* and never indicating that it should never be checked.

# Comments start with a sharp sign.	
# perl syntax requires that "." be escaped	
# Default is equivalent to ending the file with " "	
Default	2d
file: "	0
http://www.yahoo.com/.	7d
http://www.research.att.com/.	0
http://".att.com/.	1h
http://home.mcom.com/home/whatsnew/-	12h
whatsnew.html	
http://www.ncsa.uiuc.edu/SDG/Software/-	12h
Mosaic/Docs/whats-new.html	
http://snapple.cs.washington.edu:600/-	1d
mobile/	
# rarely modified	
http://www.cs.duke.edu/~pk/-	7d
HomePage.html	
# this is in my hotlist but will be different every day	
http://www.unitedmedia.com/-	never
comics/dilbert/	

Table 1: An example of the thresholds specified to *w3newer*.

2.3 Cache Consistency Issues

Determining when HTTP pages have changed is analogous to caching a file in a distributed file system and determining when the file has been modified. While file systems such as the Andrew File System [6] and Sprite [8] provide guarantees of cache consistency by issuing call-backs to hosts with invalid copies, HTTP access is closer to the traditional NFS [12] approach, in which clients check back with servers periodically for each file they access. Netscape can be configured to check the modification date of a cached page each time it is visited, once each session, or not at all. Caching servers check when a client forces a full reload, or after a time-to-live value expires.

Here the problem is complicated by the target environment: one wishes to know not only when a currently viewed page has changed, but also when a page that has not been seen in a while has changed. Fortunately, unlike with file systems, HTTP data can usually tolerate some inconsis-

tency. In the case of pages that are of interest to a user but have not been seen recently, finding out within some reasonable period of time, such as a day or a week, will usually suffice. Even if servers had a mechanism to notify all interested parties when a page has changed, immediate notification might not be worth the overhead.

Instead, one could envision using something like the Harvest replication and caching services [1] to notify interested parties in a lazy fashion. A user who expresses an interest in a page, or a browser that is currently caching a page, could register an interest in the page with its local caching service. The caching service would in turn register an interest with an Internet-wide, distributed service that would make a best effort to notify the caching service of changes in a timely fashion. (This service could potentially archive versions of HTTP pages as well.) Pages would already be replicated, with server load distributed, and the mechanism for discovering when a page changes could be left to a negotiation between the distributed repository and the content provider: either the content provider notifies the repository of changes, or the repository polls it periodically. Either way, there would not be a large number of clients polling each interesting HTTP server. Moving intelligence about HTTP caching to the server has been proposed by Gwenzman and Seltzer [3] and others.

One could also envision integrating the functionality of NO HANDS into file systems. Tools that can take actions when arbitrary files change are not widely available, though they do exist [11]. Users might like to have a unified report of new files and *W³* pages, and *w3newer* supports the "file" specification and can find out if a local file has changed. However, *snapshot* has no way to access a file on the user's (remote) file system. Moving functionality into the browser would allow individual users to take snapshots of files that are not already under the control of a versioning system such as the Revision Control System (RCS) [14]; this might be an appropriate use of a browser with client-side execution, such as *HotJava* [13].

2.4 Error Conditions

When a periodic task checks the status of a large number of URLs, a number of things can go wrong. Local problems such as network connectivity or the status of a proxy-caching server can cause all HTTP requests to fail. Proxy-caching servers are sometimes overloaded to the point of timing out large numbers of requests, and a background task that retrieves many URLs in a short time can aggravate their condition. *W3newer* should therefore be able to detect cases when it should abort and try again later (preferably in time for the user to see an updated report).

At the same time, a number of errors can arise with indi-

Tracking and Viewing Changes on the Web (NO HANDS)

vidual URLs. They can move, with or without leaving a forwarding pointer. The server for a URL can be deactivated or renamed. They may disallow retrieval by "robots," meaning that any program that follows the "robot exclusion protocol" [10] will not retrieve them. Since the cost of retrieving modification dates is small in comparison to the cost of retrieving robots.txt (part of the exclusion protocol), it may well be appropriate to ignore the robot exclusion protocol for this task, or to check robots.txt only occasionally on each host. Observing the protocol will still be advisable for hosts on which many URLs are checked, especially if the pages' contents are retrieved each time.

Finally, automatic detection of modifications based on information such as modification date and checksum can lead to the generation of "junk mail" as "noisy" modifications trigger change notifications. For instance, pages that report the number of times they have been accessed, or embed the current time, will look different every time they are retrieved.

W3newer attempts to address these issues by the following steps:

- If a URL is inaccessible to robots, that fact is cached so the page is not accessed again unless a special flag is set when the script is invoked.
- Another flag can tell *w3newer* to treat error conditions as a successful check as far as the URL's timestamp goes. For instance, if *w3newer* runs daily and checks a particular URL every four days, normally an error accessing the page on Monday will cause it to be checked again on Tuesday. With this flag, it would be checked again on Friday. In general, it seems that errors are likely to be transient, and checking the next time *w3newer* is run would be reasonable.
- When a URL is inaccessible, an error message appears in the status report, so the user can take action to remove a URL that no longer exists or repeatedly hits errors.

In addition, *w3newer* could be modified to keep a running counter of the number of times an error is encountered for a particular URL, or to skip subsequent URLs for a host if a host or network error (such as "timeout" or "network unreachable") has already occurred. Addressing the problem of "noisy" modifications will require heuristics to examine the differences at a semantic level.

3 Snapshots: External Representations of Version Histories

In addition to providing a mechanism for determining when *W³* pages have been modified, there must be a way to ac-

cess multiple versions of a page for the purposes of comparison. This section describes methods for maintaining version histories and several issues that arise with our solution.

3.1 Alternative Approaches

There are three possible approaches for providing versioning of *W³* pages: making each content provider keep a history of all versions, making each user keep this history, or storing the version histories on an external server.

Server-side support Each server could store a history of its pages and provide a mechanism to use that history to produce marked-up pages that highlight changes. This method requires arbitrary content providers to provide versioning and differencing, so it is not practical, although it is desirable to support this feature when the content provider is willing. (See Section 6.1.)

Client-side support Each user could run a program that would store items in the hotlist locally, and run *htmldiff* against a locally saved copy. This method requires that every page of interest be saved by every user, which is unattractive as the number of pages in the average user's hotlist increases, and it also requires the ability to run *htmldiff* on every platform that runs a *W³* browser. Storing the pages referenced by the hotlist may not be too unreasonable, since programs like Netscape may cache pages locally anyway. There are other external tools such as *warmlist* [16] that provide this functionality.

External service Our approach is to run a service that is separate from both the content provider and the client. Pages can be registered with the service via an HTML form, and differences can be retrieved in the same fashion. Once a page is stored with the service, subsequent requests to remember the state of the page result in an RCS "check-in" operation that saves only the differences between the page and its previously checked-in version. Thus, except for pages that change in many respects at once, the storage overhead is minimal beyond the need to save a copy of the page in the first place.

Drawbacks to the "external service" approach are that the service must remember the state of every page that anyone who uses the service has indicated an interest in, and must know which user has seen which version of each page. The first issue is primarily one of resource allocation, and is not expected to be a significant issue unless the service is used by a great many clients on a number of large pages. The second issue is addressed by using RCS's support for timestamps and requesting a page as it existed at a particular

Tracking and Viewing Changes on the Web (NO HANDS)

time. Alternatively, a version number could be retained for each <user,URL> combination.

Relative links become a problem when a page is moved away from the machine that originally provided it. If the source were passed along unmodified, then the W^3 browser would consider links to be relative to the CGI directory containing the *snapshot* script. HTML supports a BASE directive that makes relative links relative to a different URL, which mostly addresses this problem; however, Netscape 1.1N treats internal links within such a document to be relative to the new BASE as well, which can cause the browser to jump between the *htmldiff* output and the original document unexpectedly.

3.2 System Issues

The *snapshot* facility must address four important issues: use of CGI, synchronization, resource utilization, and security/privacy.

CGI is a problem because there is no way for *snapshot* to interact with the user and the user's browser, other than by sending HTML output. When a CGI script is invoked, *httpd* sets up a default timeout, and if the script does not generate output for a full timeout interval, *httpd* will return an error to the browser. This was a problem for *snapshot* because the script might have to retrieve a page over the Internet and then do a time-consuming comparison against an archived version. The server does not tell *snapshot* what a reasonable timeout interval might be for any subsequent retrievals; instead this is hard-coded into the script. In order to keep the HTTP connection alive, *snapshot* forks a child process that generates one space character (ignored by the W^3 browser) every several seconds while the parent is retrieving a page or executing *htmldiff*.

Synchronization between simultaneous users of the facility is complicated by the use of multiple files for bookkeeping. The system must synchronize access to the RCS repository, the locally cached copy of the HTML document, and the control files that record which version of each page a user has seen. Currently this is done by using UNIX file locking on both a per-URL lock file and the per-user control file. Ideally the locks could be queued such that if multiple users request the same page simultaneously, the second *snapshot* process would just wait for the page and then return, rather than repeating the work. This is not so important for making snapshots, in which case a proxy-caching server can respond to the second request quickly and RCS can easily determine that nothing has changed, but there is no reason to run *htmldiff* twice on the same data.

The latter point relates to the general issue of resource utilization. *Snapshot* has the potential to use large amounts of both processing and disk space. The need to execute

htmldiff on the server can result in high processor loads if the facility is heavily used. These loads can be alleviated by caching the output of *htmldiff* for a while, so many users who have seen versions N and $N+1$ of a page could retrieve *htmldiff*(page $_N$, page $_{N+1}$) with a single invocation of *htmldiff*. The facility could also impose a limit on the number of simultaneous users, or replicate itself among multiple computers, as many W^3 service.

Disk space is potentially a problem if the repository can grow without bound and with no cost to its users. In fact, before a service like this could be placed on the Internet, it would have to authenticate each user and limit the user to a fixed number of URLs and/or disk blocks. Most likely, one would use an Internet commerce facility to charge a fee in exchange for permission to store a collection of URLs; this fee could easily offset the cost of the storage medium since it would also be paying for the differencing service.

Lastly, security and privacy are important. Because the CGI scripts run with minimal privileges, from an account to which many people have access, the data in the repository is vulnerable to any CGI script and any user with access to the CGI area. Data in this repository can be browsed, altered, or deleted. In order to use the facility one must give an identifier (currently one's email address, which anyone can specify) that is used subsequently to compare version numbers. Browsing the repository can therefore indicate which user has an interest in which page, how often the user has saved a new checkpoint, and so on.

By moving to an authenticated system on a secure machine, one could break some of these connections and obscure individuals' activities while providing better security. The repository would associate impersonal account identifiers with a set of URLs and version numbers, and passwords would be needed to access one of these accounts. Whoever administers this facility, however, will still have information about which user accesses which pages, unless the account creation can be done anonymously.

4 Htmldiff: Comparison of HTML pages

In our experience, only a small fraction of pages on the W^3 contain information that allows users to ascertain how the pages have changed—examples include icons that highlight recent additions, a link to a "change log", or a special "What's New" page. As was mentioned in the introduction, these approaches suffer from deficiencies. They are intended to be viewed by all users, but users will visit the pages at different intervals and have different ideas of "what's new". In addition, the maintainer must explicitly generate the list of recent changes, usually by manually

Tracking and Viewing Changes on the Web (NO HANDS)

marking up the HTML.

Automatic comparison of HTML pages and generation of marked-up pages frees the HTML provider from having to determine what's new and creating new or modified HTML pages to point to the differences. There are many ways to compare documents and many ways to present the results. This section describes various models for the comparison of HTML documents, our comparison algorithm, and issues involved in presenting the results of the comparison.

4.1 What's in a Diff?

HTML separates content (raw text) from markups. While many markups (such as `<P>`, `<I>`, and `<HR>`) simply change the formatting and presentation of the raw text, certain markups such as images (``) and hypertext references (``) are "content-defining." Whitespace in a document does not provide any content (except perhaps inside a `<PRE>`), and should not impact comparison.

At one extreme, one can view an HTML document as merely a sequence of words and "content-defining" markups. Markups that are not "content-defining" as well as whitespace are ignored for the purposes of comparison. The fact that the text inside `<P>...</P>` is logically grouped together as a paragraph is lost. As a result, if one took the text of a paragraph comprised of four sentences and turned it into a list (`...`) of four sentences (each starting with ``), no difference would be flagged because the content matches exactly.

At the other extreme, one can view HTML as a hierarchical document and compare the parse tree or abstract syntax tree representations of the documents, using subtree equality (or some weaker measure) as a basis for comparison. In this case, a subtree representing a paragraph (`<P>...</P>`) might be incomparable with a subtree representing a list (`...`). The example of replacing a paragraph with a list would be flagged as both a content and format change.

We view an HTML document as a sequence of sentences and "sentence-breaking" markups (such as `<P>`, `<HR>`, ``, or `<H1>`) where a "sentence" is a sequence of words and certain (non-sentence-breaking) markups (such as `` or `<A>`). A "sentence" contains at most one English sentence, but may be a fragment of an English sentence. All markups are represented and are compared, regardless of whether or not those markups are "content-defining." In the paragraph-to-list example, the comparison would show no change to content, but a change to the formatting.

We apply Hirschberg's solution to the longest common

subsequence (LCS) problem [4, 5] (with several speed optimizations) to compare HTML documents. This is the well-known comparison algorithm used by the Unix *diff* utility [7]. The LCS problem is to find a (not necessarily contiguous) common subsequence of two sequences of tokens that has the longest length (or greatest weight). Tokens not in the LCS represent changes. In Unix *diff*, a token is a textual line and each line has weight equal to 1. In *htmldiff*, a token is either a sentence-breaking markup or a sentence, which consists of a sequence of words and non-sentence-breaking markups. Note that the definition of sentence is not recursive: sentences cannot contain sentences. A simple lexical analysis of an HTML document creates the token sequence and converts the case of the markup name and associated (variable,value) pairs to upper-case; parsing is not required.

We now describe how the weighted LCS algorithm compares two tokens and computes a non-negative weight reflecting the degree to which they match (a weight of 0 denotes no match). Sentence-breaking markups can only match sentence-breaking markups. They must be identical (modulo whitespace, case, and reordering of (variable,value) pairs) in order to match (see section 3.3 for a discussion of the ramifications of this). A match has weight equal to 1. Sentences can match only sentences, but sentences need not be identical to match one another. We use two steps to determine whether or not two sentences match. The first step uses sentence length as a comparison metric. Sentence length is defined to be the number of words and "content-defining" markups such as `` or `<A>` in a sentence. Markups such as `` or `<I>` are not counted. If the lengths of two sentences are not "sufficiently close," then they do not match. Otherwise, the second step computes the LCS of the two sentences (where words matching exactly against words are assigned weight 1, and markups match exactly against markups, as before). Let W be the number of words and content-defining markups in the LCS of the two sentences and let L be the sum of the lengths of the two sentences. If the percentage $(2 * W) / L$ is sufficiently large, then the sentences match with weight W . Otherwise, they do not match.

4.2 Presentation of the differences

The comparison algorithm outlined above yields a mapping from the tokens of the old document to the tokens of the new document. Tokens that have a mapping are termed "common"; tokens that are in the old (new) document but have no counterpart in the new (old) are "old" ("new"). We refer to the "old" and "new" tokens as "differences".

We investigated three basic ways to present the differences by creating HTML documents that highlight the dif-

Tracking and Viewing Changes on the Web (NO HANDS)

ferences with a variety of markup techniques:

Side-by-Side A side-by-side presentation of the documents with common text vertically synchronized is a very popular and pleasing way to display the differences between documents (see, for example, Unix *sdiff* or SGI's graphical diff tool *gdiff*). Unfortunately, there is no good mechanism in place with current HTML and browser technology that allows such synchronization (although it might be possible to make a document that contained a table with a document per column in which rows of the table were used to achieve synchronization).

Only Differences Show only differences (old and new) and eliminate the common part (as done in Unix *diff*). This optimizes for the "common" case, where there is much in common between the documents. This is especially useful for very large documents but can be confusing because of the loss of surrounding common context. Another problem with this approach is that an HTML document comprised of an interleaving of old and new fragments might be syntactically incorrect.

Merged-page Create an HTML page that summarizes all of the common, new, and old material. This has the advantage that the common material is displayed just once (unlike the side-by-side presentation). However, incorporating two pages into one again raises the danger of creating syntactically or semantically incorrect HTML (consider converting a list of items into a table, for example).

Our preference is to present the differences in the merged-page format to provide context and use internal hypertext references to link the differences together in a chain so the user can quickly jump from difference to difference. We currently deal with the syntactic/semantic problem of merging by eliminating all old markups from the merged page (note that this doesn't mean all markups in the older document, just the ones classified as "old" by the comparison algorithm). As a result, old hypertext references and images do not appear in the merged page (of course, since they were deleted they may not be accessible anyway). However, by reversing the sense of "old" and "new" one can create a merged page with the old markups intact and the new deleted. A more Draconian option would be to leave out all old material. In this case, there are no syntactic problems given that the most recent page is syntactically correct to begin with: the merged page is simply the most recent page plus some markups to point to the new material. We are exploring other ways to create a merged page.

An example of *hmlldiff*'s merged-page output appears in Figure 1. Markups are used to highlight old and new material as follows. Two small arrow images are used to point

to areas in the document that have changed. A red arrow points to old content and a green arrow points to new content. The arrows are also internal hypertext references to one another, linked in a chain to allow quick traversal of the differences. A banner at the front of the document contains a link to the first difference. Old text is displayed in "struck-out" font using `<STRIKE>`, which in our experience is rarely used in HTML found on the Web. Unfortunately, there is no ideal font for showing "new" text. We currently use `<I>`. Ideally, we would like to be able to color code the text or text background to highlight old and new text, but this capability is not provided by current browsers. Another approach would be to choose a font that is not active at the point of the difference.

Note that not all changes in the documents are highlighted. For example, new markups that are not "content-defining" (such as `<P>`) are not marked up. However, markups such as anchors are highlighted. Consider the example of changing the URL in an anchor but not the content surrounded by `<A>...`. In this case, an arrow will point to the text of the anchor, but the text itself will be in its original font, signifying a change to just the URL.

4.3 Issues and Extensions

Since *hmlldiff* can parse an HTML document and rectify certain syntactic problems, such as mismatched or missing markups, the only real problem it is likely to encounter is a set of changes that are so pervasive as to make the resulting merged HTML unreadable. For instance, if every other line were changed, then the mixture of unrelated struck-out and emphasized text would be muddled. We are experimenting with methods for varying the degree to which old and new text can be interspersed, as well as thresholds to specify when the changes are too numerous to display meaningfully.

Currently, *hmlldiff* is neither "version-aware" nor "web-aware". That is, *hmlldiff* only compares the text of two HTML pages. It does not compare versions of the entities that the pages refer to, access them, or invoke itself recursively on other referenced pages. This has a number of consequences. The good news is that *hmlldiff* does not incur the overhead of pulling versions from a repository or sending requests over the Web for information. This cost is consumed by *w3newer* and *snapshot*. The bad news is that some differences may be ignored. For example, if the contents of an image file are changed but the URL of the file does not, then the URL in the page will not be flagged as changed. To support such comparison would require some sort of versioning of referenced entities and would also require *hmlldiff* to have access to the version repositories. Full versioning of all entities would allow interesting com-

Tracking and Viewing Changes on the Web (NO HANDS)

BLANK PAGE

Tracking and Viewing Changes on the Web (NO HANDS)

parsons to be done, but would dramatically increase storage requirements. A cheaper alternative would be to store a checksum of each entry and use the checksums to determine if something has changed. We are exploring how to efficiently perform such "smarter" comparisons.

5 Integrating the tools

There are two entry points to NO HANDS, one through *w3newer* and one through *snapshot*.

Currently, *w3newer* is invoked directly by the user, probably by a *cron* entry, and generates an HTML document indicating which pages have changed. If specified, *w3newer* will associate three links with each document in the list:

Remember Send the URL to the *snapshot* facility, to save a copy of the page. Though the page is retrieved, the RCS *ci* command ensures that it is not saved if it is unchanged, from the previous time it was stored away.

Diff Have the *snapshot* facility invoke *hmdiff* to display the changes in a page since it was last saved away by the user.

History Have *snapshot* display a full log of versions of this page, with the ability to run *hmdiff* on any pair of versions or to view a particular version directly. (See Figure 2.)

Thus, each page that is reported as "new" can immediately be passed to *hmdiff*, and any page in the list can be "remembered" for future use. An example of *w3newer*'s output appears in Figure 3.

A user may also choose to enter *snapshot* directly to check-in pages, or view the current page or the version history. Figure 4 shows the interface to NO HANDS through *snapshot*. If the user selects the *history* link, the page shown in Figure 2 is presented. Finally, selecting two pages to compare invokes *hmdiff*, as in Figure 1.

One disadvantage of the current approach is that there is no direct interaction between *w3newer*, *snapshot*, and the *W3* browser. Viewing a page with *hmdiff* does not cause the browser to record that the page has just been seen; instead, the browser records the URL that was used to invoke *hmdiff* in the first place. Subsequently, *w3newer* uses the obsolete timestamp from the browser and continues to report that the page has been modified more recently than the browser has seen it. As a result, the user must view a page directly as well as via *hmdiff* in order to both remove it from the list of modified pages and see the actual differences.

6 Extensions

This section describes some possible extensions to the work already presented. Section 6.1 discusses an interface between RCS and *hmdiff* that is already implemented, while Sections 6.2 and 6.3 presents unimplemented extensions to integrate tracking modifications into the server and to invoke scripts via the HTTP POST protocol.

6.1 Server-side Version Control

The tools described above do not require any changes to arbitrary servers or clients on the *W3*. Existing GET and POST protocols are used to communicate with specific servers that save versions of documents and provide marked-up versions showing how they have changed. However, if a server runs *hmdiff* and some *perl* scripts, it can provide a direct version-control interface and avoid the need to store copies of its HTML documents elsewhere.

The *perl* scripts we have written provide an interface to RCS [14]. A CGI script (*/cgi-bin/rlog*) converts the output of *rlog* into HTML, showing the user a history of the document with links to view any specific version or to see the differences between two versions. Another script (*/cgi-bin/co*) displays a version of a document under RCS control, while still another (*/cgi-bin/rcsdiff*) displays the differences. If the file's name ends in *.html* then *hmdiff* is used to display the differences, rather than the *rcsdiff* program.

As an example, one might set up a Last-Modified field at the bottom of an HTML document to be a link to the *rlog* script, with the document name specified as a parameter. After clicking on this unobtrusive field, the user would be able to see the history of the document.

6.2 Server-side URL Tracking

Currently, *w3newer* runs on the user's machine, so multiple instantiations of the script may perform the same work. Although it runs a related daemon on the same machine as an AT&T-wide proxy-caching server, which returns information about pages that are currently cached on the server and may eliminate some accesses over the Internet, there is insufficient locality in that cache for it to eliminate a significant fraction of requests.

Alternatively, *w3newer* could be run on the set of pages that have been saved by the *snapshot* daemon. Regardless of how many users have registered an interest in a page, it need only be checked once: if changed, the new version could be saved automatically. Then a user could request a list of all pages that have been saved away, and get an indication of which pages have changed since they were saved by the user.

Tracking and Viewing Changes on the Web (NO HANDS)

BLANK PAGE

Tracking and Viewing Changes on the Web (NO HANDS)

Adding this functionality would be useful, since it would offer economies of scale. It would have the disadvantage of being decoupled from a given user's *W³* browser history; i.e., if a user views a page directly, the *snapshot* facility would have no indication of this and might present the page as having been modified.

6.3 Interactions with CGI scripts

Because NO HANDS can handle arbitrary URLs, it can interact with CGI scripts that use the GET protocol by passing arguments to the script as part of the URL. However, services that use POST cannot be accessed, because the input to the services is not stored.

Both *w3newer* and *snapshot* would have to be modified to support the POST protocol, in order to invoke a service and see if the result has changed, and then to store away the result and display the changes if it has. The interface to NO HANDS to support POST is unclear, however. A user could manually save the source to an HTML form and change the URL the form invokes to be something provided by NO HANDS. It, in turn, would have to make a copy of its input to pass along to the actual service. The result would be an HTTP equivalent of a UNIX *pipe*, interposing an extra service between the browser and the service the user is trying

to invoke.

Instead, the browser could be modified to have better support for forms:

- It should store the filled-out version of a form in its bookmark file, so the user could jump directly to the output of a CGI script.
- It should be able to pass a form directly to NO HANDS, along with the URL specified in the FORM tag, so that the output could be stored under RCS.

7 Conclusions

NO HANDS combines notification, archiving, and differencing of *W³* pages into a single cohesive tool. It achieves economies of scale by avoiding unnecessary HTTP accesses, saving pages at most once each time they are modified (regardless of the number of users who track it), and using RCS as the underlying versioning system. Automatic generation of differences within the HTML framework provides users with the ability to see both insertions and deletions in a convenient fashion.

In the general setting of the *W³* and document retrieval, NO HANDS benefits two communities: users of the *W³* and

Tracking and Viewing Changes on the Web (NO HANDS)

longer have to browse to find pages of interest that have changed: HTML providers no longer have to create suitably marked-up pages to show "what's new". While such automation is clearly helpful in this general context, we expect that NO HANDS will be a critical part of more focused uses of the W^3 , especially in areas involving collaborative and distributed work.

Several issues still need to be addressed. In particular, many of the complications of NO HANDS could be avoided by better integration with W^3 browsers and servers. For instance, viewing the difference between an older version of a page and its current version should update the browser's notion of when the page was last visited. Finally, the increasing availability of distributed, hierarchical HTTP repositories such as Harvest [1] will be both an opportunity and a challenge for scalable notification mechanisms and version archives.

- [7] J. W. Hunt and M. D. McIlroy. An algorithm for differential file comparison. Technical Report Computing Science TR #41, Bell Laboratories, Murray Hill N.J., 1975.
- [8] M. Nelson, B. Welch, and J. Ousterhout. Caching in the Sprite network file system. *ACM Transactions on Computer Systems*, 6(1):134-154, February 1988.
- [9] M. Newbery, Katipo, <http://www.vuw.ac.nz/newbery/Katipo.html>
- [10] A standard for robot exclusion, <http://web.nexor.co.uk/mak/doc/robots/norobots.html>
- [11] David S. Rosenblum and Balachander Krishnamurthy. Generalized event-action handling. In Balachander Krishnamurthy, editor, *Practical Reusable UNIX Software*, chapter 9. John Wiley & Sons, New York, 1995.
- [12] R. Sandberg, D. Goldberg, S. Kleiman, D. Walsh, and B. Lyon. Design and implementation of the Sun network filesystem. In *Proceedings of the USENIX 1985 Summer Conference*, pages 119-130, June 1985.
- [13] Sun Microsystems. *The HotJava Browser. A White Paper*. Available as <http://java.sun.com/1.0alpha3/doc/overview/hotjava/browser-whitepaper.ps>.
- [14] W. Tichy. RCS: a system for version control. *Software—Practice & Experience*, 15(7):637-654, July 1985.
- [15] Url-minder, <http://www.neumind.com/URL-minder/URL-minder.html>
- [16] Warmlist, <http://glimpse.cs.arizona.edu:1994/paul/warmlist>

References

- [1] C. Mic Bowman et al. Harvest: A scalable, customizable discovery and access system. Technical Report CU-CS-732-94, Dept. of Computer Science, University of Colorado-Boulder, March 1995.
- [2] B. B. Cutler III. w3new, <http://www.suff.com/bcutler/programs/w3new/w3new.html>.
- [3] James S. Gwertzman and Margo Seltzer. The case for geographical push-caching. In *Proceedings of the Fifth Workshop in Hot Topics in Operating Systems (HOTOS-V)*, pages 51-55, Orcas Island, WA, May 1995. IEEE.
- [4] D. S. Hirschberg. A linear space algorithm for computing maximal common subsequences. *Communications of the ACM*, 18(6):341-343, June 1975.
- [5] D. S. Hirschberg. Algorithms for the longest common subsequence problem. *Journal of the ACM*, 24(4):664-675, October 1977.
- [6] J. Howard et al. Scale and performance in a distributed file system. *ACM Transactions on Computer Systems*, 6(1):51-81, February 1988.

WO 97/15890

PCT/US96/17142

COMPUTER PROGRAM LISTING

[illegible]

```

let
fun gen_list () { _ _ _ x
}
| gen_list (x::xs) { _ _ _ x
if new_only then
}
else
{OLD,s1::gen_list xs | 0 0 {} }
} gen_list [] {y::yes _-[] =
(MEM,y)::gen_list [] yn 0 {} }
| gen_list [] p2 p3 d1 as {x,r,d1::xs} =
if p1 = n andalso p2 = y then
(case d1 of
OLDMEM(d1) =>
let
val arrive,d1 = dr
val MEMM(d1) = had 12)
in
gen_list b1 b2 0 0 d1
end
| h => {SUCCEED})
if p1 = x then
else if p1 = z then
(MEM,had 12)::gen_list [] {} 123 p1 (p2+1) d1)
else
{if new_only then
}
else
{OLDM,had 11}::
gen_list [] b1 12 (p1+1) p3 d1))
end
fun elim_rebuild_h_list =
let
(* combine around MEMBASE *)
fun f_ellm5 ((MEM,s)::MEMBASE,y)::(MEMM,s)::xs =
f_ellm5 ((MEM,SQO((x,r,d1)::xs))
| f_ellm5 (xs::xs) = m::f_ellm5 xs)
| f_ellm5 [] = [])
(* get rid of MEMs around terms containing
only WACURS *)
fun f_ellm6 ((MEMM,s)::xs) =
if is_variable s then MEMM else MEMBASE
s::f_ellm6 xs)
| f_ellm6 (xs::xs) = x::f_ellm6 xs))
| f_ellm6 [] = []
(* collapse alternation of OLD/MEM/OLD/MEM *)
fun f_ellm7 ((OLD,e)::(MEMM,s1)::
{OLDM,M2::(MEMM,M2)::xs} =
f_ellm7 ({OLDM,SQO([e],o2)}::xs)
| f_ellm7 ({OLDM,SQO([e],o2)}::xs)
| f_ellm7 (xs::xs) = x::f_ellm7 xs))
| f_ellm7 [] = []
(* eliminate sequences of MEMs (OLMs) *)
fun f_ellm8 (((t1,s1)::((t2,p1)::xs) =
if t1 == SUCCEED andalso p1 = 12 then

```

[illegible]

45

```

fun fill at st
  fill_for esc_bool esc_for end fun
  let
    fun iter i =
      if i < length!at then
        update!(at,i,fill_for i);
        if (esc_bool i) then
          esc_for i;
        else
          iter (i+1);
        end
      else
        end_for i;
      end
    in
      iter at
    end
  end

fun sample (fill) at at fill_for fun =
  fill at at fill_for
  (fun x => false) (fun x => () (fun s => ()
    fun comp_weight (col,row) (w,v) if i =
      let
        val diffiv_1 =
          (weight
            (vector.sub(iv_1,i),
              vector.sub(iv_1,j fill)
            )
          in
            max(sub(iv_1,i,w),
              max(sub(iv_1,i-1),sub(iv_1,i)))
          end
        fun verify_colm (col,row) (x,y) =
          if col <= s then
            fill_for row
            comp_weight
              (col,row) (ref 0, ref 0 col)
            (fun x => (sub(Col,row) > sub(row,x)))
            (fun s =>
              (sample(fill) Col (row)
                (fun s => (sub(Col,row))));
              loop_ref(Col,row));
              verify_row (col,y) (col,row)
            (fun s =>
              loop_ref(Col,row);
              verify_colm (col,row) (x,y))
            else
              (x,y)
            end
          end
        and verify_row (col,row) (x,y) =
          if row <= y then
            fill_for col
            comp_weight
              (row,col) (ref 0, ref 0 row)
            (fun s => (sub(row,s) > sub(row,s)))
            (fun s =>
              (sample(fill) row (row)
                (fun s => (sub(row,s))));
              loop_ref(row,col));
              verify_col (col,y) (col,row)
            (fun s =>
              loop_ref(row,col);
              verify_row (col,row) (x,y))
            else
              (x,y)
            end
          end
        end
      end
    end
  end
end

```



```

if (piv a <= w then
  ip_diff (i-1) (j)
  Vector.tabulate(ipiv,a-1),
  (for a <>
    Vector.sub(A,(i-1-1)))
  Vector.sub(ipiv,y-1),
  (for y >
    Vector.sub(B,(y-1-1)))
end

else
  DIFF(w,(i))
  in
    ident_prefix (ipiv,x1) (piv,y1)
    (letw,(revres) real)
  end
end

else
  DIFF(wt,(rev real))
  and ident_prefix a y (wt,real) =
  if (Debug) then
    (print "prefix "; print a;
     print " "; print y;
     print "a") else ();
    if a <= w and also y <= a then
      let wt DIFF(w,1) =
        weight (Vector.sub(A,a-1),
                Vector.sub(B,y-1))
      in
        if w <= 0 then
          (ident_prefix (x1) (y-1))
          (letw,(a-1,y-1,DIFF(w,1))
           :real)
        else
          main a y (wt,real)
        end
      else
        DIFF(wt,(rev real))
      end
    end
  in
    main 1 1 (0,())
  end
end

end

structure Chain : WEIGHT =
struct
  type T = char
  fun weight diff (a:T,b:T) =
    if a = b then DIFF(1,()) else DIFF(0,())
  end
end

structure CharDiff = DIFF(Chain,structure W = Chain);

```

Oct 23, 1995
17:18:06

/home/tball/sml/mlweb/mlweb.sml

king — 09:32 Sep 18

1

4

```

signature NAME =
sig
  type HTMLatt
  datatype HTMLDoc =
    | MARKUP of string * (HTMLatt list)
    | DOCUMENT of string * (HTMLatt list)
    | PROLOG of string * (HTMLatt list)
    | EPILOG of (HTMLDoc list)
    | SEQ of (HTMLDoc list)

  (* get the names of input fields - TEXTAREA, INPUT, SELECT *)
  val info : string -> string option list * bool * bool * bool

  val info_stream : instream * string -> string option list * bool * bool * bool

  val parse : string -> HTMLDoc

  val parse1 : string -> HTMLDoc

  val parse_stream : instream * string -> HTMLDoc

  val unparse : HTMLDoc -> string

  val unparse_pa : (HTMLDoc -> bool) -> (HTMLDoc -> string) -> HTMLDoc -> string

  val dump : (HTMLDoc -> HTMLDoc) -> HTMLDoc -> HTMLDoc

  val getatt : string -> (string * string) list -> string option

  val hier : HTMLDoc -> HTMLDoc

  val expl : HTMLDoc -> HTMLDoc list

  val sentences : HTMLDoc -> HTMLDoc list
end

structure HTML : HTML =
struct
  open StringUtil
  exception BadHtml
  type HTMLatt = string * string
  datatype HTMLDoc =
    | MARKUP of string * (HTMLatt list)
    | DOCUMENT of string * (HTMLatt list)
    | PROLOG of string * (HTMLatt list)
    | EPILOG of string * (HTMLatt list)
    | SEQ of (HTMLDoc list)

  val toUpper : JcrLib.Locater
  fun maybe_str s : string =
    if Index.find s (".", ";") = NONE then expandStr s ^ ""
    else handle Notfound s
  fun unparse_attr (in_v::_) : string =
    if in_v = "" then unparse_attr (in_v::_)
    else unparse_attr (in_v::_)
  fun unparse_attr nil : string =
    unparse_attr (in_v::_)
  fun unparse (HTML s) : string =
    unparse (SEQ s)
  fun unparse (DOCUMENT s) : string =
    unparse (SEQ s)
  fun unparse (PROLOG s) : string =
    unparse (SEQ s)
  fun unparse (EPILOG s) : string =
    unparse (SEQ s)
  fun unparse_pa pat act (HTML s) : string =
    if pat (HTML s) then (act (HTML s)) else s
  fun unparse_pa pat act (SEQ s) : string =
    if pat (SEQ s) then (act (SEQ s))
    else (fold op "" (map unparse_pa pat act) s)
  fun unparse_pa pat act (DOCUMENT s) : string =
    if pat (DOCUMENT s) then (act (DOCUMENT s))
    else (fold op "" (map unparse_pa pat act) s)
  fun unparse_pa pat act (PROLOG s) : string =
    if pat (PROLOG s) then (act (PROLOG s))
    else (fold op "" (map unparse_attr s))
  fun unparse_pa pat act (EPILOG s) : string =
    if pat (EPILOG s) then (act (EPILOG s))
    else (fold op "" (map unparse_attr s))
  fun unparse_pa pat act (HTML s) : string =
    if pat (HTML s) then (act (HTML s))
    else (fold op "" (map unparse_attr s))
  fun unparse_pa pat act (SEQ s) : string =
    if pat (SEQ s) then (act (SEQ s))
    else (fold op "" (map unparse_attr s))
  fun unparse_pa pat act (DOCUMENT s) : string =
    if pat (DOCUMENT s) then (act (DOCUMENT s))
    else (fold op "" (map unparse_attr s))
  fun unparse_pa pat act (PROLOG s) : string =
    if pat (PROLOG s) then (act (PROLOG s))
    else (fold op "" (map unparse_attr s))
  fun unparse_pa pat act (EPILOG s) : string =
    if pat (EPILOG s) then (act (EPILOG s))
    else (fold op "" (map unparse_attr s))
  fun explode_ravs : string -> string list
  fun expl (HTML s) : HTMLDoc list =
    let

```


4

4

/home/tball/sml/mlweb/mlweb.sml

king — 09:32 Sep 18

Oct 23, 1995
17:18:06

```

let
  fun ffields' (h::t) =
    if fields' t then ffields h
    | fields' [] = []
  in
    ffields' hmls
  end
  { fields' = []
  fun submit (markup, attlist) =
    case # of
      'HTML' => [
        case getatt "type" attlist of
          SONG => ffields (temp s = "submit"
          | NONE => false)
        | _ => false
        ]
      | submit (SQ(hmls)) =
    let
      fun submit' (h::t) =
        if submit h then true
        else submit' t
      | submit' [] = false
    in
      submit' hmls
    end
    | submit _ = false
    fun has_markup desired (markup(s, _)) =
      if ffields (temp s) = desired
      | has_markup desired (SQ(hmls)) =
    let
      fun has_markup' (h::t) =
        if has_markup desired h then true
        else has_markup' t
      | has_markup' [] = false
    in
      has_markup' hmls
    end
    | has_markup _ = false
  } ffields hml, submit hml, has_markup 'BODY' hml, has_markup '/BODY'
in
  hml
end
fun info s =
let
  in
    info_stream (info, s) before close_in info
  end
end

```

Oct 23, 1995
17:18:15

ome/tball/sml/mlweb/htr. ex

king — 17:19 Jul 25

1

1

```

datatype lexresult =
  OPEN |
  ENDOPEN |
  DECLOPEN |
  PROCOPEN |
  CLOSE |
  OPENCOMMENT |
  CLOSECOMMENT |
  EQ |
  CONTENT of string |
  COMMENT of string |
  STR of string |
  ID of string |
  EOF

val eof = fn () => EOF
val filename = ref ""
(*
fun dprint(s) = (print (s ^ "(" ^ (makestring (0)) ^ ")\n"))
*)
fun dprint(s) = ()

val error = fn (x,l:int) => output(std_err,x ^ " line " ^
  (makestring (l+1)) ^
  ", file '" ^ (!filename) ^ "'\n")

**
** MU COMM COMMSTUP;
**structure HtmlLex
**  ws      = ([\ \t \n])~;
**  nonws   = ([^ \ \t \n <=])~;
**  count =
**
<INITIAL>["<>"] => (dprint "CONTENT"; CONTENT(yytext));
<INITIAL>"<" => (YYBEGIN MU; dprint "OPEN"; OPEN);
<INITIAL>">" => (YYBEGIN MU; dprint "ENDOPEN"; ENDOPEN);
<INITIAL>"!" => (YYBEGIN MU; dprint "DECLOPEN"; DECLOPEN);
<INITIAL>"?" => (YYBEGIN MU; dprint "PROCOPE"; PROCOPEN);
<MU>">" => (YYBEGIN INITIAL; dprint "CLOSE"; CLOSE);
<MU>"--" => (YYBEGIN COMM; dprint "OPENCOMMENT"; OPENCOMMENT);
<MU>"=" => (dprint "EQ"; EQ);
<MU>"\"{^\"}\n\"}\" => (dprint "STR"; STR(substring(yytext, 1, String.size(yytext)-2)));
<MU>(nonws) => (dprint "ID"; ID(yytext));
<MU>(ws) => (lex());
<COMM>"--" => (YYBEGIN MU; dprint "CLOSECOMMENT"; CLOSECOMMENT);
<COMM>"<" => (YYBEGIN INITIAL; dprint "CLOSECOMMENT";
(yybufpos := (yybufpos)-1); CLOSECOMMENT);
<COMM>["<>"] => (dprint "COMMENTTEXT"; COMMENT(yytext));
      => (error("html lex: ignoring bad character " ^ yytext, !yylineno); lex());
**

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king -- 10:28 Sep 12

1
17

```

#!/usr/local/bin/perl -- # A comment mentioning perl, to prevent looping
eval "exec perl -s $0 $@" if $? != 0;
if 0;
# For the latest info on Fred Douglas's update of v1.0 of w3newer, see
# http://www1126.research.att.com/~douglas/track_urls/
# For the latest information on Brooks Cutler's version of w3new, check the URL
# http://www.eng.paradyme.com/cgi-bin/bbcburn/user/bcutterspag-w3new
# w3newer v1.0 - Creates a What's New list of http URL's
# Modified by Fred Douglas <douglas@research.att.com>
# to separate user interface from polling.
# w3newer is a program that will extract a list of URL's from either your mosaic
# hotlist, or will extract the URL's from a HTML document. It will then obtain
# the HTTP/1.0 modification dates for each document listed, and output a HTML
# file with the URL's sorted by their last modification time
#
# Package: w3newer
# Author: Brooks Cutler (bcutler@paradyme.com)
# Author: Fred Douglas (douglas@research.att.com)
# Latest version: 1.0
# Last updated: March, 1995
# Archive: w3newer.tar.gz
# (includes everything you need to run it, except perl v4.036)
#
# What it does.
#
# 1. w3newer first extracts the URL's from your hotlist or a HTML document, and
#    cached modification dates and access times from its own cached
#    history and the WWW browser's history file.
# 2. w3newer then hands off the URL's to a separate process, which
#    proceeds to do a HTTP/1.0 HEAD on each http: URL that might
#    be out of date and stores the last modified time of each URL (if
#    available).
# 3. Finally, w3newer sorts its output based on the last modified time of the
#    document and categorizes them by the month they were modified. For
#    non-http URL's or URL's for which it can't retrieve the last modified
#    time, it will list these at the bottom of its HTML output.
#
# This program was written because Brooks found himself frequently checking web
# pages in his hotlist to see if they had recently been changed. Now
# he runs w3newer from cron nightly, and checks its output each morning.
#
# pointing w3newer at a list of URL's
# =====
#
# 1. Using your Mosaic hotlist by default, w3newer will read your Mosaic
#    hotlist file and extract the document URL's and document titles. It will
#    look for your hotlist in your home directory as the file
#    ~/.mosaic-hotlist-default
#
#    this can be overridden by using the environment variable
#    W3NEWER_HOTLIST
#
#    The default and environment variable can be overridden with the
#    command line arguments -hotlist_fn or -f with a argument of the file in
#    Mosaic-hotlist format.
# 2. Extracting URL's from a HTML document if you call w3newer with the
#    -u argument and a HTML document URL, it will retrieve that document,
#    extract all the links in the document, and check the modification time of
#    each link.
# 3. Extracting URL's from a HTML document and inline modification times
#    if you call w3newer with the -html argument and a HTML document
#    URL, it will extract the document links and retrieve the modification
#    times, but unlike #2 (-u argument), it will replace the modification
#    times within the document rather than producing a sorted list.
#
# When w3newer parses the document, it will look for a tag of the form
#
# <w3newer url="URL">
#
# for example
# <w3newer url="http://www.next.com/file.html">
#
# when it finds a tag like the one above, it will retrieve the last

```

Oct 24, 1995
09:12:02/home/tball/src/w3newer/w3newer
king - 10:28 Sep 12

2

17

```

modification time for the document specified by the url- line. And
include the last modification of that URL (if it exists)

for more information on the usage, run w3newer with the -h argument

The latest information on this program can be found on the web - see above

This program was built with software written by others.
Their contributions are greatly appreciated..

libwww-perl v0.60 or later by Roy Fielding (fielding@cs.ucl.ac.uk)
http://www.ics.ucl.ac.uk/webSoft/libwww-perl/
evaluate_parameters (evap) by Stephen C. Lidie (slidie@lehigh.edu)
All the implementations of evaluate_parameters are available via
anonymous FTP from ftp.lehigh.edu (128.180.6) 4). Look in the
pub/evap/evap-2.x directory for the latest compressed tar file.
perl v4.036 by Larry Wall (lwall@netlabs.com)
ftpable from ftp.uu.net in /systems/gnu as perl-4.036.tar.gz
w3new v0.4 by Brooks Cutter (bcutter@paradyme.com)

known bugs
*****

o none (that I am aware of...)

Future work
*****

o none (probably). I started to write some code that would do a checksum
on a document specified by a URL ... this would allow y. to track the
modification time based on when the content changed... however I
never finished coding/testing the checksum piece, and abandoned it.
you'll find the code commented out below..

o picked up 1/95 by P. Douglas. Cached modification times are positive;
cached checksums are negative.

o eventually I hope to convert this into a WWW applet and let the
browser do the work. You'll find more info on WWW Applets at
http://www.let.rug.nl/~bert/W3A/W3A.html

o Other future work (FD): multi-user "database" of info. Types other than
http.

Brooks Cutter (bcutter@paradyme.com)
http://www.sng.paradyme.com/cgi-bin/bbcurn?user=bcutter

-----
w3newer is distributed under the Artistic License (included with your Perl
distribution files and with the distribution of the libwww-perl package).
-----

To run this program either set the variable below with the location
$ENV{'W3NEWER'} = "/home/tball/src/w3newer" unless ($ENV{'W3NEWER'});
$ENV{'W3NEWER'} = "/home/tball/src/w3newer";
$w3newer_dir = $ENV{'W3NEWER'};
unless($w3newer_dir) {
    print "EOF";
}

***ERROR*** Unable to start w3newer

Sorry, I don't know where w3newer support files are installed!

Please either set the variable $w3newer_dir near the top of the w3newer script,
or set the W3NEWER environment variable to the base package directory.

EOF

exit: 0

$process_url = "perl $w3newer_dir/lib/w3newer-est.pl";

unless($ENV{'http_proxy'}) {
    $ENV{'no_proxy'} = "www126.research.att.com.sno.att.com.radiash.research.att.com.localhost.tbw.att.com.no.att.com";
    $ENV{'http_proxy'} = "http://www.tb.att.com:8080/";
}

o What version are you using? (set this if using $w3newer_dir)
$libwww_ver = "0.60";
$ENV{'LIBWWW_PERL'} = "$w3newer_dir/lib/libwww-perl-$libwww_ver";
o If defined $ENV{'LIBWWW_PERL'}
$ENV{'LIBWWW_PERL'} = "$w3newer_dir";

```


24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

3

17

```

#!/usr/bin/perl
# where's libwww-perl?
# Either in the directory specified by LIBWWW_PEARL
# or under lib in $w3newer_dir
$ENV{'LIBWWW_PEARL'} =
"/home/douglib/lib/perl".

$w3newer_dir = "perl $ENV{'LIBWWW_PEARL'}/POST http://www.tbu.att.com/~douglib/cgi-bin/w3newer.cgi";

# architecture-specific library (for sys/socket.ph)
require "arch.pl";
$arch = $arch;
if (1-d ($libloc = "$w3newer_dir/lib/arch/$arch")) {
    unshift(@INC, $libloc);
}

# TRAIL: Let's keep some usage stats
open(MAIL, ">maillog.tball");
$date = "date";
print MAIL "Subject: w3newer\n$ENV{'USER'}\n$date\n";
close(MAIL);
# TRAIL

require "time.pl";
require "request.pl"; # part of www nov
require "stat.pl";
require "www.pl";
require "wwwurl.pl";
require "wwwerror.pl";
require "wwwdata.pl";
require "wwwbot.pl"; # part of libwww-perl v0.20 or later..

unshift(@INC, "/appl/evap/lib") if 1-d "/appl/evap";
require "evap.pl";

require "parse_html.pl";

# You shouldn't need to change anything below here..
$pl_request_timeout = 45;

@months = (
    "January", "February", "March", "April", "May", "June", "July",
    "August", "September", "October", "November", "December"
);

@days = ("Sunday", "Monday", "Tuesday", "Wednesday", "Thursday", "Friday", "Saturday");

$w3newer_author_home = "stuff.com";
$w3newer_ver = "1.0";
$w3newer_author_email = "bcutter@w3newer_author_home";
$w3newer_author = "<EDF>";
<a href="http://www.w3newer_author_home/cgi-bin/bbcoutn?user=bcutter">
Brooks Cutter</a>
EDF
$w3newer_oldurl =
"http://www.stuff.com/~bcutter/home/programs/w3new/w3new.html";
$w3newer_url =
"http://www126.research.att.com/~douglib/track_urls/";

$w3newer_author2 = "<EDF>";
<A HREF="http://www.research.att.com/orgs/est/people/douglib/">Fred Douglas</A>
EDF

# if URL isn't fully qualified, this is used with $wwwurl::absolute;
$base = "file://localhost: ($ENV{'PWD'}) [ ($ENV{'cwd'}) ] .";

$UserAgent = "w3newer/$w3newer_ver ($wwwLibrary)"; # Set up User-Agent: header
# Set the default User-Agent
$www::set_def_header("http://User-Agent:$UserAgent");

$hotlist_netscape = ("($ENV{'HOME'})/.netscape-bookmarks.html",
    "$ENV{'HOME'})/.NCOM-bookmarks.html");
$history_netscape = ("($ENV{'HOME'})/.netscape-history",
    "$ENV{'HOME'})/.NCOM-global-history");
$hotlist_mosaic = ("($ENV{'HOME'})/.mosaic-hotlist-default",
    "$ENV{'HOME'})/.mosaic-global-history");
$history_mosaic = ("($ENV{'HOME'})/.mosaic-global-history");

$default_browser = "netscape";

# This requires the evap.pl library..
$evap::set::EDF;
}

w3newer

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

4

17

```

# the following options control output processing
verbose, v: switch
output_fn, o: file = W3NEWER_MTHL, "~/public_html/w3newer_output.html"
# if you prefer to have it output to STDOUT, uncomment the next line..
# output_fn, o: file = W3NEWER_MTHL, stdout
# the following control where to get URLs
mosaic, m: switch
metacape, m: switch
hotlist_fn, h: file = W3NEWER_HOTLIST, ""
history_fn, h: file = W3NEWER_HISTORY, ""
url, u: string
html_url, hml, string
# these control processing
from_address, from: string = EMAIL_ADDRESS, ""
checksum, cs: boolean = TRUE
# these control the format of the output
display_header, dh: boolean = TRUE
display_footer, df: boolean = TRUE
display_unknown, du: boolean = TRUE
display_recent, dr: boolean = TRUE
display_unchanged, du: boolean = TRUE
embedden, b: switch
last_visited, lv: switch
timestamp_source, ts: switch
# these control cached information
cache_modules, cm: boolean = TRUE
cm_file, cmf: file = W3NEWER_CMFILE, "~/w3newer_modules"
cm_def_threshold, cmdt: string = W3NEWER_CM_DEF_THRESHOLD, "10"
want_latest, vl: boolean = FALSE
cm_obsolete_threshold, cmot: string = W3NEWER_CM_OBS_THRESHOLD, "10"
cache_thresholds_fn, cfn: file = W3NEWER_CTFILE, "~/w3newer_thresholds"
# miscellaneous options, esp. for debugging
debug, d: integer = 0
max_urls, max_urls: integer = -1
multi_user, mu: switch
nop, n: switch = false
# integrated into modtimes
w3newer_cs_file, csf: file = W3NEWER_CSFILE, "~/w3newer_checksums"
snapshot, s: switch
snapshot_url, su: string = W3NEWER_SNAP_URL, "http://radish.research.att.com/cgi-bin/no_bands"
ignore_snapots, ign: switch
format, f: switch
FORMAT optional_file_list
EOF
sevap_pdt = split("/n/sevap_pdt");
sevap_md += EOF;
w3newer vivw3newer_viv

```

w3newer is a program that will extract a list of URL's from a hotlist (either Metacape or Mosaic format), or will extract the URL's from an HTML document. It checks the modification dates of URLs that have not been checked recently and produces a report of URLs that have been modified since they were viewed.

- It tries not to poll things too often. It uses the following techniques:
- If your browser's history shows something has been accessed recently, it isn't checked.
 - If w3newer checked something recently, it isn't checked.
 - If w3newer knows something has changed since it was last visited, it marks it as new but doesn't check for a more current date unless it deems the modification date information to be "obsolete" (see documentation for cm_obsolete_threshold).
 - The definition of "recent" can vary from URL to URL, based on specifications in a "thresholds" file.

Currently, w3newer only handles http URLs, not ftp, file, or others. This may change.

Examples

```

w3newer
w3newer -usage_help
w3newer -v
w3newer -full_help
w3newer -ns -m -lv -v (metacape, bold, last visited, verbose)
verbose

```

When verbose is set, w3newer will print out the name of each URL it checks and generate some other status information.

When checksum is set, w3newer will use checksums to tell if something has been modified if the modification date is unavailable.

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

5

17

output_fn
This specifies where the HTML output is placed (stdout indicates standard output).

display_header
determines whether or not to print a header on the HTML output

display_footer
determines whether or not to print a footer on the HTML output

display_unknown
determines whether or not to display hosts for which it cannot find HTTP/1.0 last modification dates

display_recent
determines whether or not to display URLs that were recently accessed and not checked during this run.

display_unchanged
determines whether or not to display URLs that have not changed since they were last seen.

hotlist_fn
The filename of your Netscape/Mosaic hotlist. By default it is set to the file `--netscape-bookmarks.html` or `--mosaic-hotlist-default`, depending on whether the "netscape" or "mosaic" option is set. The default can be overridden with the environment variable `W3NEWER_HOTLIST`.

history_fn
The filename of your Netscape/Mosaic history. By default it is set to the file `--netscape-global-history` or `--mosaic-global-history`, depending on whether the "netscape" or "mosaic" option is set. The default can be overridden with the environment variable `W3NEWER_HISTORY`.

url
`W3newer` will retrieve the HTML file specified by this argument, extract all hyperlinks and images from the document, and proceed to check the last modification time for each URL.

html_url
If you call `W3newer` with the `-u` argument and a HTML document URL, it will extract the document links and retrieve the modification times, but unlike the previous option `-u` argument, it will replace the modification times within the document rather than producing a sorted list.

When `W3newer` parses the document, it will look for a tag of the form

```
<W3newer url="URL">
```

For example

```
<W3newer url="http://www.host.dom/file.html">
```

When it finds a tag like the one above, it will retrieve the last modification time for the document specified by the `url` line, and include the last modification of that URL if it exists.

from_address
By default, the program will set the HTTP/1.0 From: header to your Username (specified by `USER` environment variable) and your hostname. If your hostname isn't fully qualified (ie: doesn't have a "." in it) then it will run `domainname` and tack on the result (if any).

This value is also used for the snapshot arguments (see below).

This value can be overridden by setting the environment variable `EMAIL_ADDRESS` or passing your email address with the `-from` flag.

multi_user
This flag is not yet implemented. It refers to sharing modification dates among multiple users.

netscape
This flag indicates that Netscape bookmark and history files should be used.

mosaic
This flag indicates that Mosaic hotlist and history files should be used.

emodifier
This flag indicates that certain items that are deemed to be of special interest should be generated in boldface. Currently this refers to the dates of items that have been modified since they were last seen.

last_visited
This flag indicates that the time when a URL was last visited, according to its history file, should be included in `W3newer` output.

timestamp_source
This flag indicates that the source of information for a URL (proxy:reaching server, polling, previous run, etc.) should be included in `W3newer` output.

mtime_fn
This option specifies the file that carries modification dates of URLs. The default can be overridden by the environment variable

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king - 10:28 Sep 12

6
17

```

#NEWER_CACHE.
cache_modtimes
This option indicates that modification dates should be cached. If
FALSE, cm_file is ignored.
cm_def_threshold
This option specifies the default frequency for checking URLs. in
the form mm, hh, dd, ww for M minutes, hours, days, or weeks, respectively
want_latest
If this switch is set, always get the current modification date of a
URL that has not been accessed recently, rather than skipping it if
its known modification date is more recent than the time it was
visited by the browser and not beyond the date specified by the
-cmnot switch
cm_absolute_threshold
This switch specifies when a cached modification date should not be
considered to be current even if it indicates that a URL has not
been seen since it was modified.
cache_thresholds_in
This option specifies the file that specifies how often different
URLs should be checked. If it does not exist, then all URLs are
checked with a frequency determined by the -cmnot option. The format
of this file is:
URL whitespace time_specification
Default: time_specification
Onsolete: time_specification
Default: overrides the value of -cmnot. Onsolete overrides -cmnot. If
time_specification is "never" the URL is skipped even if it is found
in the hotlist. The URL can be a pattern to match against perl regexp
syntax.
debug
When debug is > 0, additional debugging output is generated. The higher
the number, the more output there is.
max_urls
This option specifies a limit on the number of URLs to process,
primarily for debugging purposes.
snapshot
If set, the output will contain links to the snapshot CGI daemon.
snapshot_url
URL to use for snapshot.
more_hosts
Set in order to bypass saved information about which URLs disallow robots (all) or none.
forms
Set to produce a form for each URL rather than direct hyperlinks.
EOF
evap_cm = split(/\\/, $evap_cm);
do {
  Error returned from evap\n if (! $evap($evap_pdt, $evap_cm) == 1);
  if from_address was set, set the default header..
  if ($options{'from_address'}) {
    www$set_def_header('http', 'From', $options{'from_address'});
  } else {
    local($user) = $ENV{'USER'} || getlogin || (getpwnam($c)){'0'} || 'Intruder!!!';
    local($host) = $ENV{'HOST'} || chopisfoo = '/usr/ucb/hostname' || 'Unknown_host';
    local($host) = $hostname.PQDN;
    $options{'from_address'} = join(' ', $user, $hostname.PQDN);
  }
  $option_string = '';
  for ("Diff", "Banner", "History") {
    $option_string .= "options $_";
  }
  if ($options{'output_fn'} eq '-' || ($options{'output_fn'} =~ /^s$/)) {
    open(OUT, $options{'output_fn'})
    || die "Can't write to $options{'output_fn'}: $!\n";
  } else {
    local($oldmask) = umask;
    if (!-e $options{'output_fn'}) {
      local($lib) = stat(1);
      local($mode) = $ary[EST_MODE];
      umask($oldmask | (-mode & 0777));
      if (!-e $options{'output_fn'}) {
        rename($options{'output_fn'}, "$options{'output_fn'}-") ||
          print STDERR "Warning: can't rename $options{'output_fn'}\n";
      }
    }
    print(OUT, $options{'output_fn'})
    || die "Can't write to $options{'output_fn'}: $!\n";
    mask $oldmask;
  }
  $options{'verbose'} = 1 if ($options{'debug'}) & debug implies verbose.
}

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

7

17

```

if (!options('message') || !options('netcape')) {
    options('hotlist_in') || options('url') || options('html_url') || {
        options('default_browser') = 'TRUE';
    }
}

sub newchurl {
    local ($url) = @_;
    if (options('debug') > 0) {
        print STDERR "Get a URL ($url) we didn't ask for. We know about:";
        local ($u);
        for $u (sort keys %urls) {
            print STDERR "  $u\n";
        }
    }
}

sub find_first_existing {
    local ($names) = @_;
    local ($file);
    foreach $file ($names) {
        if (-e $file) {
            return $file;
        }
    }
    print STDERR "Can't find an existing file from: " . join(" ", $names);
    if (options('debug')) {
        return "";
    }
}

if (options('message')) {
    options('hotlist_in') = find_first_existing(@hotlist_message);
    options('history_in') = find_first_existing(@history_message);
} elsif (options('netcape')) {
    options('hotlist_in') = find_first_existing(@hotlist_netcape);
    options('history_in') = find_first_existing(@history_netcape);
}

load_checksums if options('checksum');

# read in cached modification dates & checksums, if any. Then if we
# see URLs with cached info, update the cached info based on any new
# info, e.g. "last seen" dates.

if (options('cache_modtimes')) {
    %cached_modtimes = read_cached_modtimes(options('cache_file'));
} else {
    %cached_modtimes = ();
}

if (options('history_in') ne "") {
    %visit_times = load_history(options('history_in'));
}

if (options('url')) {
    $url = %wwwurl{absolute($base, options('url'))};
    if (!%wwwnot_allowed{$url, $UserAgent}) {0}
        die
        "Sorry, but the remote site does not allow robots to retrieve that URL.\n";
    %urls = %extract{linux_desc{$url}};
} elsif (options('html_url')) {
    $url = %wwwurl{absolute($base, options('html_url'))};
    if (!%wwwnot_allowed{$url, $UserAgent}) {0}
        die
        "Sorry, but the remote site does not allow robots to retrieve that URL.\n";
}

print OUT <display.html>{$url};
exit;
} else {
    # use 'hotlist_*' names to have qualified pattern matching on call by ref.
    # is this necessary?
    local(%hotlist_urls, %hotlist_dates);
    load_hotlist(options('hotlist_in'), %hotlist_urls, %hotlist_dates);
    %urls = %hotlist_urls;
    %urls{accessed} = %hotlist_dates;
    for $url (keys %urls) {
        if (%urls{accessed}{$url} < %visit_times{$url}) {
            print STDERR "Warning: Visit time from hotlist for $url lower than visited history. " if (options('debug') && %visit_times{$url} > 0;
            %visit_times{$url} = %urls{accessed}{$url};
        }
    }
}

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king -- 10:28 Sep 12

8

17

```

    $timestamp = "notinit";

    if (length $opts{times} == 0) {
        delete $options{autoload};
    }

    sub threshold {
        local ($url) = @_;
        return threshold($url) if defined(threshold($url));
        local ($u);
        foreach $u (@thresholds) {
            if ($u =~ m/^$url$/i) {
                print STDERR "Using threshold for %s for URL %s\n" $u, $url if $options{debug} > 0;
                return threshold($u);
            }
        }
        return $options{cache_def_threshold};
    }

    # initialize units in terms of seconds
    local($units) = "m 60. s 60. d 24. w 7.";
    local($current) = 1;
    local($seconds, $seconds1, $seconds2);
    while ($seconds > 0) {
        ($seconds1, $seconds2, $seconds) = $units;
        $seconds1/$seconds2 = $seconds1 * $current;
        $seconds1 = $seconds1/$seconds2;
    }

    sub str2secs {
        local ($orig_thresh) = @_;
        local ($m, $sec, $min, $st) = 0;
        $st = $orig_thresh;
        while ($st ne "") {
            if ($st =~ /never/) {
                return -1;
            }
            if ($st =~ /(\d+)([a-z])$/) {
                $m = $1; $sec = $2; $st = $3;
                $m = $m * $seconds1/$seconds2;
                $m = $m * $seconds1/$seconds2;
            }
            if ($st =~ /(\d+)([a-z])$/) {
                $m = $1; $sec = $2;
                return $m;
            }
            if ($st =~ /(\d+)([a-z])$/) {
                $m = $1; $sec = $2;
                return $m;
            }
            die "Can't parse threshold $orig_thresh";
        }
        return $m;
    }

    sub sec2str {
        local ($s) = @_;
        local ($m, $sec, $min, $st) = 0;
        if ($s < 0) {
            warn "Negative time ($s) passed to routine";
            return "";
        }
        $s = $s;
        for $u ("m", "d", "h", "w") {
            die "Units not initialized" if !defined($units{$u});
            $m = 0;
            while ($s >= $units{$u}) {
                $m++;
                $s = $s - $units{$u};
            }
            $res = "$m:$u" if $m;
        }
        $res = "$s:s" if $s;
        $res = "" if $res eq "";
        return $res;
    }

    if ($options{cache_thresholds}) {
        local $default_cache_threshold = 0;
        warn "WARN: $options{cache_thresholds} [1] die 'Can't open $options{cache_thresholds}';
        print STDERR "Thresholds: %s" if $options{debug} > 0;
        while ($res) {
            if ($res =~ /(\d+)([a-z])$/) {
                $m = $1; $sec = $2;
            }
        }
    }

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

9

17

```

Sections["cm_def_threshold"] = &str2secs(1);
set_default = 1;
printf(STEPPR "%10s %d\n", "Default", Options["cm_def_threshold"]); if (Options["debug"] > 2;
) else if ("Obsolete" < 1) {
    Options["cm_obsolete_threshold"] = &str2secs(1);
    set_obsolete = 1;
    printf(STEPPR "%10s %d\n", "Obsolete", Options["cm_obsolete_threshold"]); if (Options["debug"] > 2;
) else if ("") {
    next; // skip comment
} else {
    local(surl, snum) = split(/%/);
    if (snum ne "") {
        threshold(surl) = &str2secs(snum);
        push(thresholds, surl);
    }
    printf(STEPPR "%10s %d\n", surl, threshold(surl)); if (Options["debug"] > 2;
)
}
close(THRESH);
if ("set_default") {
    Options["cm_def_threshold"] = &str2secs(Options["cm_def_threshold"]);
}
if ("set_obsolete") {
    Options["cm_obsolete_threshold"] =
    &str2secs(Options["cm_obsolete_threshold"]);
}

surls_processed = 0;
new_urls = &old_urls + &skipped_urls + &new_checked + &old_checked +
    &inaccessible_urls + &nohits = {};

print STEPPR "Current time is ", time, "\n" if (Options["debug"]);

# create pipes to talk to the server -- this should be changed to use POST
# directly.
pipe(GETURLS, SENDURLS); if (die "Can't create pipe: $?");
pipe(GETRESULTS, SENDRESULTS); if (die "Can't create pipe: $?");
if ($pid = fork) {
    close(GETURLS);
    close(SENDRESULTS);
} else if (defined $pid) {
    close(GETURLS);
    close(GETRESULTS);
    open(STDOUT, ">&SENDRESULTS"); if (die "Can't dup stdout");
    open(STDIN, "<&GETURLS"); if (die "Can't dup stdin");
    exec $w3newer ($?) die "Can't exec POST: $?";
} else {
    die "Can't fork: $?";
}

$? = 0;

for $url, $new $surls {
    local($t) = threshold($url);
    next if ($t == -1);
    printf STEPPR "CM_SAV: %s url threshold %d\n" if (Options["verbose"]);
    next unless($url =~ /^http:/);
    next if ($url =~ /^http:\/\//); // skip scripts
    printf(STEPPR "%sURLid=%s\n", ($? > 0) ? "s" : "", $?, $url);
    $? = 0;
}
close(GETURLS);
while(GETRESULTS) {
    printf(STEPPR "%s\n", $? if (Options["debug"] > 1;
    next if "http:";
    chop;
    ($url, $mod $when, $rest) = split(/%/);
    die "Bad input from back end" if $rest ne "";
    $surl = $url; next if (defined $surl);
    next if $mod eq "N/A";
    local($cached_mod, $cached_checked, $nohits) =
        split($?, $cached_mod($url));
    $nohits = "" if (Options["ignore_nohits"]);
    if ($mod < $cached_mod) {
        $mod = $cached_mod;
        $when = $cached_checked;
        printf(STEPPR "SP: %s url mod %d when %d\n" if (Options["verbose"]);
        $when = $cached_checked if $cached_checked < $when;
    }
}

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

to
17

```

    $cached_modtimes($url) = join(' ', $mod, $when, $nobots);
    $stamp($url) = "proxy-caching server";
} else { $mod > 0 } {
    print STDERR "old info for $url\n" if ($options{'debug'} > 0;
}

# create pipes to talk to the local back end.
pipe(GETURLS, SENDURLS) || die "Can't create pipe: $!";
pipe(GETURLS, SENDRESULTS) || die "Can't create pipe: $!";
if ($pid = fork) {
    close(GETURLS);
    close(SENDURLS);
} else { defined $pid } {
    close(SENDURLS);
    close(GETURLS);
    open(STDOUT, ">>SENDRESULTS") || die "Can't dup stdout";
    open(STDIN, "<<GETURLS") || die "Can't dup stdin";
    exec $process_urls || die "Can't exec $process_urls: $!";
} else {
    die "Can't fork: $!";
}

for $url (keys %urls) {
    local($t) = $threshold{$url};
    if ($t == -1) {
        print STDERR "Skipping ignored URL $url\n" if ($options{'debug'});
        next;
    }
    print STDERR "CHK: $url threshold $t\n" if ($options{'verbose'});
    next unless $url =~ /^http:/ || $url =~ /^file:/;
    next if ($url =~ m/^http:\/\?/); # skip scripts
    last if ($options{'max_urls'} >= 0 && --$urls_processed > $options{'max_urls'});

    $last_checked = $mod = $visit_times{$url}; $cached_cs = 0;
    if ($url =~ /^http:/) {
        if (defined $cached_modtimes{$url}) {
            ($cached_mod, $cached_checked, $nobots) =
                split(' ', $cached_modtimes{$url});
            if ($cached_checked >= $last_checked) {
                $last_checked = $cached_checked;
                $stamp($url) = "previous run";
            }
            if ($cached_mod > 0) {
                if ($cached_mod > $visit_times{$url} &&
                    (time - $cached_mod >=
                     $options{'rm_obsolete_threshold'}) &&
                    $options{'want_latest'}) {
                    push(@new_urls, $url);
                    $url_time{$url} = $cached_mod;
                    print STDERR "URL known to be new; skipping\n" if ($options{'verbose'});
                    if ($options{'debug'}) {
                        local($mdate) = localtime($url_time{$url});
                        local($vdate);
                        $vdate = localtime($visit_times{$url}) if $visit_times{$url};
                        $vdate = "Never visited" unless $visit_times{$url};

                        chop $mdate;
                        chop $vdate;
                        print STDERR "visited $vdate ($visit_times{$url}) modified $mdate\n";
                    }
                }
                $mod = $cached_mod;
            } else { $cached_mod < 0 } {
                $cached_cs = - $cached_mod;
            }
        }
        if ($nobots eq "none") {
            push(@noccs_urls, $url);
            print STDERR "URL known not to allow robots; skipping\n" if ($options{'verbose'});
            next;
        }
        print STDERR "URL $url visited $visit_times{$url} mod $mod\n" if ($options{'debug'});
        # store the max for later - used for knowing when to add emphasis
        $url_checked{$url} = $last_checked;
        if ($last_checked < $t && time) {
            print(STDERR "Using cached $url age to threshold to [$stamp($url)]\n";
            time = $last_checked; $t = 0 if ($options{'verbose'});
            if ($mod > $visit_times{$url}) {
                push(@new_urls, $url);
            }
        }
    }
}

```


Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

11

17

```

    } else {
        push(@skipped_urls, $url);
    }
    $url_time{$url} = $mod if $mod > 0;
    next;
}
print STDERR "Checking $url\n" if ($options{verbose});
next if ($options{nop});
# need to change this to handle missing files, down hosts, etc -- be
# smarter about return codes.
local($succeeded) = sprintf("%s to %s\n", $url, $last_checked,
                             ($cached_mod > 0) * $cached_mod, 0);
print STDERR $succeeded;
print STDERR "---- $succeeded" if ($options{debug}) > 0;
else {
    # file
    local($file) = $url;
    $file =~ s/^file:--//;
    $file =~ s/^file:SDV("MOK")//;
    $file =~ s/^/localhost//;
    if ( -e $file ) {
        local($ary) = stat($file);
        $time = $ary{ST_MTIME};
        $url_time{$url} = $time;
        $stamp{$url} = "local filesystem";
        print STDERR "File $file exists -- $time\n" if ($options{debug});
        if ($time > $visit_time{$url}) {
            push(@new_urls, $url);
        } else {
            push(@old_urls, $url);
        }
    } else {
        push(@inaccessible_urls, $url);
        print STDERR "File $file does not exist\n" if ($options{debug});
        next;
    }
}
}
close(STDOUT);
while (<<GETRESULTS>) {
    print STDERR "---- $_" if ($options{debug}) > 0;
    chop;
    ($url, $mod, $status, $resp, $rest) = split(/ /);
    die "Bad input from back end" if ($rest ne "");
    if ($status eq "MOK") {
        push(@new_urls, $url);
        $cached_modtimes{$url} = join($_, $url_time{$url}, $visit_time{$url}, "nobots");
        next;
    }
    if ($status eq "MOK") {
        $status{$url} = $resp if $resp ne "";
        push(@inaccessible_urls, $url);
        # don't check again until timeout
        # $url_checked{$url} = time;
        $cached_modtimes{$url} = join($_, $url_time{$url}, time);
        next;
    }
    if ($status ne "OK") {
        print STDERR "Warning: bad status $status from back end, url $url\n";
    }
    if ($mod > 0) {
        local ($changed);
        if ($mod <= $visit_time{$url}) {
            push(@old_urls, $url);
            $changed = 0;
        } else {
            push(@new_urls, $url);
            $changed = 1;
        }
        $url_time{$url} = $mod;
        $cached_modtimes{$url} = join($_, $mod, time);
        print STDERR "OK: $url ($mod) $changed\n", $changed ? "un" : "un";
        if ($options{verbose});
        $stamp{$url} = "polled";
    } else {
        if ($mod < 0) {
            print STDERR "Checksum($url):$cached_cs == $dn\n", $mod;
            if ($options{debug});
            $cached_modtimes{$url} = join($_, $new_cs, time);
            if ($cached_cs <= $new_cs) {
                push(@new_urls, $url);
            } else {

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

King - 10:28 Sep 12

12
17

```

        push($old_chksums, $url);
    }
    else {
        print STDERR "NOK: $url resp $respin" if ($options{verbose});
        push($unaccessible_urls, $url);
    }

    exit if ($options{nop});

    print OUT "<html>.\n";
    if ($options{display_header}) {
        $now($date = localtime);
        print OUT "<<EOF";
        <title>Status of URLs</title>
        <!--This page generated by <a href="$w3newer_url">w3newer v$w3newer_version</a> on $date</h1>
    }
    print OUT "<div>url_list";
    if ($options{display_footer}) {
        print OUT "<<EOF";
    }
    <hr>
    This page generated by <a href="$w3newer_url">w3newer v$w3newer_version</a>
    a program written by $w3newer_author and modified by $w3newer_author2.
    </hr>
    print OUT "</html>.\n";
    close(OUT);

    *save_checksum if ($options{checksum});

    if ($options{cache_modtimes}) {
        write_cached_modtimes($options{cache_file}, %cached_modtimes);
    }

    exit;

sub url_desc {
    local ($url) = @_;
    if ($url =~ /http/) {
        return $url;
    }
    else {
        return $url;
    }
}

sub bydesc {
    $url_desc($a) < $url_desc($b);
}

sub display_url_by_month {
    local (@url_list) = @_;
    local ($url, $old_month, $month, $month_no, $list, $, $day, $yday, $ydaynum, $hold);
    local ($vdate, $source);
    $list = "";
    $month = "";
    for $url (@url_list) {
        $old_month = $month;
        ($day, $month, $year, $ydaynum) = (gmtime($url_time($url)))[3,4,5,6];
        $month = months($month_no);
        $yday = $days($ydaynum);
        if ($month ne $old_month) {
            $list = "</div>\n" if ($old_month);
            $list = "<div>$month Year</div>\n";
            $list = "<div>\n";
        }
        $list = "<div> <div> <a href=\"$url\">";
        $list = $url_desc($url);
        $list = "</a>";
        $hold = ($options{bolden}) &&
            ($url_time($url) > $visit_time($url)) &&
            ($visit_time($url) > 0);
        if ($options{last_visited}) && $visit_time($url) > 0 {
            $now($vdate = localtime($visit_time($url)));
            $vdate = "<div>visited $vdate";
        }
        else {
            $vdate = "";
        }
        if ($options{timestamp_source}) && $stamp($url) ne "" {
            $source = "<div>based on $stamp($url)";
        }
        else {

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

13
17

```

        $source = ""
    }

    $list = "<STRONG>" if $bold;
    $list = "Sunday, $month $day"
    $list = "</STRONG>" if $bold;
    $list = $vdate;
    $list = $source;
    $list = "\n";

    $list = $snapshot($url) if $options['snapshot'];
    $list = $snapshot($url);
    $list = "</nohtml>\n";
}

$list = "</ul>\n</ul>\n";

sub snapshot {
    local ($url) = @_;
    return "" if $url =~ "http://";
    if ($options['form']) {
        return qq<
    <form name="url" type="hidden" value="$url">
    <input name="email" type="hidden" value="$options['from_address']">
    <select name="type">
        <option string
    </select>
    <input type="submit" value="invoke">
    </form>
    </font>
    } else {
        local ($a) = qq<A href="$options['snapshot_url']?email=$options['from_address']&url=$url&type=: 0">|;
        return qq< | $aNameoverNameover</A> | $aDiffDiff</A> | $aHistoryHistory</A> : 0";
    }
}

sub display_url_list {
    local($list) = "---- This page generated by w3newer $w3newer_ver by "
    . "Brooks Cutter ($w3newer_author_email) ---\n";
    if ($new_urls) {
        $list = "<h2>The following URLs have changed:</h2>\n";
        $list = $display_url_by_month(reverse sort { $url_time($a) < $url_time($b) } $new_urls);
    }
    if ($options['checksum']) {
        if ($new_checksums) {
            $list = "<h2>The following URLs may have changed (no modification time available):</h2>\n";
            foreach $url ($new_checksums) {
                $desc = $url_desc($url);
                $list = qq<li><nohtml> <a href="$url">$desc</a>\n;
                $list = $snapshot($url) if $options['snapshot'];
                $list = "</nohtml>";
            }
            $list = "</ul>\n";
        }
        if ($old_checksums) {
            $list = "<h2>The following URLs don't seem to have changed (no modification time available):</h2>\n";
            foreach $url ($old_checksums) {
                $desc = $url_desc($url);
                $list = qq<li><nohtml> <a href="$url">$desc</a>\n;
                $list = $snapshot($url) if $options['snapshot'];
                $list = "</nohtml>";
            }
            $list = "</ul>\n";
        }
    }
    if ($options['display_recent']) && $skipped_urls {
        $list = "<h2>The following URLs were checked recently and not pulled <li>\n";
        for $url (sort bydesc $skipped_urls) {
            $desc = $url_desc($url);
            $list = qq<li><nohtml> <a href="$url">$desc</a>\n;
            if ($options['timestamp_source']) && $stamp($url) ne "" {
                local ($age) = time - $url_checked($url);
                if ($age < 0) {
                    if ($options['verbose']) {
                        local($foo, $fcol, $foo2);
                        chop($foo=>ctime($age));
                        $foo2 = $foo;
                        chop($foo=>ctime($url_checked($url)));

```

Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king — 10:28 Sep 12

14
17

```

if ($?) {
    warn "Surl apparently checked in the future";
    printf(STDERR "time is checked %s\n", $fcol, $lcol);
}

- $age = 0;

local $lcr = $secs2str/$age;
it = $lcr*$age;
local $lcr = sprintf("%s %s",
    $stamp($url) eq "history" ?
    "visited" : "checked",
    $lcr);

$list = ". based on $stamp($url): $age ago";

$list = $snapshot($url) if ($options['snapshot']);
$list = "no snapshot";

$list = "no snapshot";

if ($options['display_unchanged'] eq $old_urls) {
    $list = "The following URLs have not changed:\n";
    $list = $display_url_by_month; reverse sort; $url_time/$a/$b/$c/$d/$e $old_urls;
}

if ($options['display_unknown'] eq $inaccessible_urls) {
    $list = "Unable to access the following URLs:\n";
    for $url (sort { $url($a) < $url($b) } $inaccessible_urls) {
        $desc = $url_desc($url);
        $list = qq"<li>no</li>" . $url($url) . $desc . "\n";
        if ($url_accessed($url) > 0) {
            $age = $time($url_accessed($url));
            $list = "(last visited $age)";
        }
        $list = $snapshot($url) if ($options['snapshot']);
        $list = "no snapshot";
        if (defined($status($url))) {
            local $why = $wwwerror($status($url));
            if (defined($why)) {
                $list = "<li>$why</li>";
            }
        }
    }
    $list = "no snapshot";
}

if ($options['display_unavailable'] eq $unavailable_urls) {
    $list = "The following URLs are not accessible to robots:\n";
    for $url ($unavailable_urls) {
        $desc = $url_desc($url);
        $list = qq"<li>no</li>" . $url($url) . $desc . "\n";
        $list = $snapshot($url) if ($options['snapshot']);
        $list = "no snapshot";
    }
    $list = "no snapshot";
}

return($list);
}

# a generic routine to figure out what kind of hotlist we're looking at.
# potentially useful if the hotlist is fed in through stdin or something.
sub load_hotlist {
    local $file = $hotlist_url; $hotlist_dates = 0;
    local $descs = $times;

    open(FILE, "<$file") or die "Can't open file $file for reading: $!";

    $f = <FILE>;
    if ($f =~ /^mozilla-hotlist-format-1/) {
        load_mozilla_hotlist_or_history(FILE) $hotlist_dates $hotlist_urls;
    } else {
        die "Don't know how to handle hotlist with format $f";
    }
    close(FILE);
}

sub load_mozilla_hotlist_or_history {
    local $name = $hotlist_urls; $descs = 0;
    local $f = $f; $descs = $times;
}

```


Oct 24, 1995
09:12:02

/home/tball/src/w3newer/w3newer

king - 10:28 Sep 12

17
17

```

sub read_cached_modtimes {
    local ($file) = @_;
    local ($times, $url, $mod, $checked, $snobots);
    open(CACHED,$file) || return 1;
    while (<CACHED) {
        ($url, $checked, $mod, $snobots, $rest) = split(/\\s+/);
        die "Bad input from $file" if $rest ne "";
        $times{$url} = join(' ', $mod, $checked, $snobots);
        $stamp{$url} = "previous run";
    }
    close(CACHED);
    return $times;
}

sub write_cached_modtimes {
    local ($file,$checked) = @_;
    local ($mod, $checked, $key, $snobots);
    return if !length $key || $checked eq 0;
    rename($file, "$file-") if -e $file;
    if (!open(CACHED,">$file")) { warn "Unable to write file $file"; return; }
    foreach $key ($keys_to_cache) {
        ($mod, $checked, $snobots) = split(' ', $checked{$key});
        print CACHED "$checked\\t$mod\\t$snobots\\n";
    }
    close(CACHED);
}

```

```

# =====
# localdatetotime() is a modification of the get_gmtimefrom wwwdates.pl
#
# Description:
#   Translate a date string to machine time (seconds since Epoch)
#
# Usage:
#   $mtime = localdatetotime($date);
#
#   where $date is in the following format:
#
#   "Thu Feb 3 17:03:55 1994"
#
sub localdatetotime {
    local($_) = @_;
    local($l) = 0;
    local($l) = 0;
    local($day, $m, $yr, $hr, $min, $sec, $date, $atime, $mon, $midx);
    local($offset) = 0;

    # Split date string
    local(@v) = split;

    ($weekday, $m, $day, $atime, $yr) = @v;
    $yr += 1900;
    ($hr, $min, $sec) = split(/:/, $atime);

    if (!$mon) { ($yr += 70); return 0; }

    # Translate month name to number
    $midx = index($Mntr, substr($m,0,3));

    if ($midx < 0) { return 0; }
    else { $mon = $midx + 1; }

    # Translate to seconds since Epoch
    return (timegm($sec, $midx, $hr, $day, $mon, $yr));
}

```

TCP

Oct 23, 1995
17:19:54

/home/tball/douglis-tmp/nohandsBE
kmg — 16:16 Oct 23

1
11

```
#!/usr/local/bin/perl

# This script is a backend to the snapshot script in my CGI directory
# It is separated primarily for purposes of debugging: it doesn't rely
# on having a CGI-environment set up and instead works just with
# command-line args.

# One of the things being done for simplicity here is to just keep
# appending to the log of <url.timestamp> tuples, instead of
# overwriting in place.

# to do:
#   locking

$debug=1;
$need_locks = 1;
$maxlength = 1024 * 512;      # 512 Kbytes
$snapshot = "http://www126.research.att.com/cgi-bin/no_hands";
# this should be PWD but that may not be normalized to /home.
$snapdir = "/home/douglis/tmp/cgi/test/snapshot";

local($oldfh, $select(STDERR);
$| = 1;
select(STDOUT); $| = 1;
select($oldfh);

unless($INC, "/home/douglis/lib/perl");
# architecture-specific library (for sys/socket.pm)
require "arch.pl";
$sarch = $arch;
if (-d $libloc = "/home/douglis/arch/$sarch/lib/perl") {
    unshift($INC, $libloc);
}

($STD_CLOCKID, $P_SETFD) = (1, 2); # from solaris sys/fcntl.h

$lib_www_ver = "0.40";
unshift($INC,
    $ENV{"LIBWWW_PERL"} || "/home/douglis/lib/perl/libwww-perl-$lib_www_ver");
$:
$root = $ENV{"PWD"};

if (defined $ENV{"PATH"}) {
    # print STDERR "Path: $ENV{"PATH"}<br>" if $debug;
} else {
    $ENV{"PATH"} = "/home/douglis/arch/$sarch/bin:/home/douglis/bin:/bin:/usr/sbin:/usr/bin:/usr/ucb";
}

unless($ENV{"http_proxy"}) {
    # $ENV{"no_proxy"}="att.com";
    # want research.att.com to go outside
    $ENV{"no_proxy"}="tba.att.com,ncr.com,radiash.research.att.com,www126.research.att.com,ib.att.com,bo.att.com,cc.att.com,nc.att.com,dsi.att.com";
    $ENV{"http_proxy"}="http://radiash.research.att.com:8080/";
    # $ENV{"http_proxy"}="http://bluel.research.att.com:8080/";
}

require "cgi-lib.pl";
$EXIT_FN = "cleanup_tmp";
require "cgi-exit.pl";
require "cgi-alarm.pl";
require "duplexpipe.pl";
require "normalize_html.pl";
require "ctime.pl";
require "stat.pl";
require "www.pl";
require "lock.pl" if $need_locks;
require "flush.pl";
require "printid.pl";

if ($debug) {
    $argstring = join(" ", @ARGV);
    print STDERR $argstring;
}
```


Oct 23, 1995
17:19:54

/home/tball/douglls-tmp/nohandsBE

king - 16:16 Oct 23

2
11

```

requires "findprog.pl";
%$rcsdiff,$htmldiff,$filter_anchors%=<findprog>rcsdiff, "htmldiff", "filter_anchors";
%$rcsdiff,$htmldiff,$rlog2html%=<findprog>rcsdiff, "htmldiff", "rlog2html";
%exit_gracefully>Unable to locate rcsdiff, 0) if !defined $rcsdiff;
%exit_gracefully>Unable to locate htmldiff, 0) if !defined $htmldiff;
%exit_gracefully>Unable to locate rlog2html, 0) if !defined $rlog2html;
%exit_gracefully>Unable to locate filter_anchors, 0)
% if !defined $filter_anchors;
$rcsdir = $rcsdiff;
$rcsdir =~ s/\/rcsdiff$//;

$cs="<rcsdir/>";
$cc="<rcsdir/>";
$ccs="<rcsdir/>";

$no_hands_url =
  "http://www126.research.att.com/~douglls/no_hands/";

$no_hands_author = <<EOF:
<A HREF="http://www126.research.att.com/~douglls">Fred Douglass</A>
EOF
%no($no_hands_author);
$htmldiff_author = <<EOF:
<A HREF="http://www-spr.th.att.com/~tball">Tom Ball</A>
EOF
%no($htmldiff_author);

% exit_gracefully>Wrong number of arguments, 1) if $ARGV != 2;

% need full path here because of possible BASE directive
$logo = qq:<a href="http://www126.research.att.com/~douglls/no_hands/"
</A>\n;
$trailer = <<EOF:
</BODY>
</>
<FONT SIZE=3>
This page generated by <a href="$no_hands_url">NO HANDS</a>.
an <strong>ATT&#177; proprietary</strong> program written by $no_hands_author.
<br>
Any <a href="mailto:nonhands@tucan.research.att.com">comments</a> are appreciated.
$logo
</HTML>
EOF
$htmldiff_trailer = qq: using <a href="http://www-spr.th.att.com/~tball/htmldiff.html">htmldiff</a></> (also proprietary, <br>
<!-- written by $htmldiff_author);

($operation, $user, $url, $serverversion) = @ARGV;

@tmpfiles = ();

% ascc: array of operations and required fields. u->URL, e->email(user).
% write access, p->people
% allowable_operations = {Remember, "uue",
  "Diff", "uue",
  "Status", "ue",
  "List_All", "e",
  "History", "ue",
  "Users", "p",
  "Remember_if_new", "uue",
  "All_URLs", "p",
  "URL_Users", "up",
  "View_Current", "ue"
};

%exit_gracefully>Illegal operation $operation specified, 1)
unless defined $allowable_operations{$operation};

% print STDERR "allowable_operations($operation)=$allowable_operations{$operation}\n" if $debug;

if ($allowable_operations{$operation}) == /e/ {
  %exit_gracefully>Illegal user name, 1) if $user =~ /\w0.1-5/;
  $user = "people/$user";

```

Oct 23, 1995
17:19:54

/home/tball/douglls-tmp/nohandsBE

kdg - 16:16 Oct 23

3

11

```

if ($allowable_operations($operation) == /u/) {
    $exit_gracefully(sprintf("Illegal URL: '%s'", $url), 1);
    if ($url =~ m!^http://[\w\.-]+\.[\w-]+!$!) {
        # temp until add logging
        if ($operation eq "View_Current") {
            print "Location: $url\n\n";
            exit;
        }
        $need_rcs = 1;
    } else {
        $need_rcs = 0;
    }

    if ($allowable_operations($operation) == /p/) {
        opendir(PEOPLE, "people") || $exit_gracefully("Unable to open directory 'people'", 0);
        local($people) = readdir(PEOPLE);
        $rcs = grep(/^[a-z]+$/ $people);
        print STDERR "Checking users $rcs\n";
        join(' ', $rcs);
        if ($debug > 1) {
            close(PEOPLE);
        }

        if ($need_rcs) {
            $in_url = $url2fn($url);
            $rcsfile = $rcsfile.$in_url;

            if (-e $rcsfile) {
                printf STDERR "As previously checked in\n";
                $url: if ($debug > 0) {
                    $exit_gracefully("File $rcsfile not writable", 0) unless -w $rcsfile;
                }
                $existed = 1;
            } else {
                $existed = 0;
            }
        }

        $users = ();
        if ($operation eq "URL_users") {
            if ($operation eq "Diff" && $userversion != /v\d+/) ||
                $operation eq "List_All" ||
                $operation eq "History" || $operation == /.users/i ||
                $operation eq "Remember_if_new" || $operation eq "All_URLs") {
                foreach $times ($times) {
                    if ($operation ne "URL_users") {
                        undef $lastseen;
                        undef $lastuser;
                    }
                    if (-r $times) {
                        $exit_gracefully("file $root/$times not writable", 0) if ! -w $times;
                        open(TIMES, "<$times") || $exit_gracefully("Unable to open $root/$times: $?", 0);
                        unless (($need_locks || $lock(TIMES, 1)) {
                            $exit_gracefully("Unable to acquire lock on $times",
                                "A file was temporarily unavailable. Please try again later. (You can reload this page !)".
                        )
                        while (<TIMES) {
                            next if /"/;
                            local($sturl, $when) = split(' ', $.);
                            if ($operation eq "Diff" || $operation eq "History" ||
                                $operation eq "Remember_if_new" || $operation eq "URL_users") {
                                if ($sturl eq $url) {
                                    if ($operation eq "URL_users") {
                                        local($user) = $times;
                                        $user =~ s!^people/!!;
                                        $lastseen($user) = $when;
                                    } else {
                                        $prevseen = $lastseen if $userversion eq "previous:mate";
                                        $lastseen = $when;
                                        $lastseen = join(' ', $when);
                                    }
                                }
                            }
                        }
                    }
                    if ($operation eq "List_All" || $operation eq "Users") {
                        if ($operation eq "All_URLs") {
                            $lastseen($sturl) = $when;
                            $lastseen($sturl) = join(' ', $when);
                        }
                        if ($operation eq "All_URLs") {
                            # not here: use if set just lx
                            $lastuser($sturl) = $times;
                            $lastuser($sturl) = join(' ', $times);
                        }
                    }
                }
            }
        }
    }
}

```

Oct 23, 1995
17:19:54

/home/tball/douglis-tmp/nohandsBE

king - 16:16 Oct 23

4

11

```

slastseen = $prevseen if defined $prevseen && $userversion eq "penultimate";
close (TIMES);
$unlock(TIMES) if $need_locks;
if ($operation eq "Users") {
    local ($keys) = keys $lastseen;
    local ($user) = $times;
    $user -- s/"people"/";";
    $count($user) = $keys + 1;
    elsif ($operation eq "All_URLs") {
        for ($keys $lastseen) {
            $count{$_}++;
        }
    }
    elsif ($operation eq "Users") {
        print "Nothing has ever been checked in.";
        print $trailer;
        exit 0;
    }
    if ($operation eq "List_All") {
        print $title("Pages recorded for $user");
        print "<pre>Sorted by Domain/</pre>\n";
        print $logo;
        print "<TABLE> <TR> <TH>URL</TH><TH>Modification timestamp</TH></TR>\n";
        for (sort bydomain keys $lastseen) {
            local ($url, $rest) = split(' ', $lastseen{$url});
            printf qq:<TR><TD><A HREF="%s?url=%s&email=%s&type=History">%s</A></TD><TD>%s</TD></TR>\n, $snapshot, $url, $user, $count{$url};
        }
        print "</TABLE>". $trailer;
        exit 0;
    }
    elsif ($operation eq "Users") {
        print $title("Users of the NO HANDS service");
        print "<TABLE> <TR> <TH>User</TH><TH>Pages recorded</TH></TR>\n";
        for (reverse sort byusercount keys $count) {
            printf qq:<TR><TD><A HREF="%s?email=%s&type=ListAll">%s</A></TD><TD>%d</TD></TR>\n, $snapshot, $_, $count{$_};
        }
        print "</TABLE>". $trailer;
        exit 0;
    }
    elsif ($operation eq "URL_users") {
        print $title("Users who have recorded $url");
        print $logo;
        for (sort keys $lastseen) {
            printf qq:<A HREF="%s?email=%s&type=ListAll">%s</A> (%s<BR>)\n, $snapshot, $_, $lastseen{$url};
        }
        print $trailer;
        exit 0;
    }
    elsif ($operation eq "All_URLs") {
        print $title("Pages recorded by the NO HANDS service");
        print $logo;
        print "<TABLE> <TR> <TH>URL</TH><TH>User(s)</TH><TH>Last version</TH></TR>\n";
        for (sort byusercountURL keys $count) {
            printf qq:<TR><TD><A HREF="%s?url=%s&email=%s&type=ViewAll20Current">%s</A></TD><TD></TD></TR>\n, $snapshot, $_, $user, $count{$url};
            if ($count{$url} > 1) {
                printf qq:<TD><A HREF="%s?url=%s&type=URLAll20Users&email=%s">users</A>.\n, $count{$url}, $snapshot, $_, $user;
            }
            else {
                printf qq:<A HREF="%s?type=ListAll20All&email=%s">%s</A>.\n, $snapshot, $lastuser{$url}, $lastuser{$url};
            }
        }
        print "</TD>";
        local ($fn) = Arc4file($url2fn($url));
        getversion: if ( -e $fn ) {
            open(RCSFILE, $fn) || warn "Can't open $fn, last getversion:
            # do we need to lock just to read 1st line?
            local ($line) = <RCSFILE>;
            if ($line =~ /^head: (\d+)\. (\d+)/) {
                printf qq:<TD><A HREF="%s?url=%s&email=%s&type=History">%s</A></TD></TR>\n, $snapshot, $_, $user, $line;
            }
        }
        close(RCSFILE);
    }
    print "</TR>\n";
    print "</TABLE>\n$trailer";
    exit 0;
    elsif ($operation eq "Remember_if_new") {

```

Oct 23, 1995
17:19:54/home/tball/douglls-tmp/nohandsBE
king — 18:18 Oct 235
11

```

    if defined($lastseen) {
        $msg = <<EOF>>
This <A href="$url">page</A> was previously saved.
<A href="$snapshoturl=$url&mail=$user&type=History">View the history</A>
EOF
        print $msg.$ trailer;
        exit 0;
    } else if defined($lastseen) && $operation eq "Diff" {
        $msg = <<EOF>>
This URL has not previously been saved by you.
<UL>
<LI> <A href="$snapshoturl=$url&mail=$user&type=Remember">Take a snapshot</A>
<LI> <A href="$url">View current version</A>
<LI> <A href="$snapshoturl=$url&mail=$user&type=History">View the history</A>
</UL>
EOF
        $exit_gracefully("User tried to access a URL not previously seen".
            $msg);
    } else {
        $exit_gracefully("The RCS history associated with the url <a href='$url'>$url</a> is unavailable". $);
    }

    local $headers,$content,$response;

    if ($operation eq "Diff" || $operation eq "/Remember.") || $operation eq "History" {
        $current = "pages/$in_url";
        if !-e $current {
            unlink($current); $exit_gracefully("unable to unlink $current". 0);
        }
        if ($need_locks) {
            $lockfile = "locks/$in_url";
            open(LOCK, ">$lockfile" || $exit_gracefully("Can't lock $lockfile". 0);
            push(@tmpfiles, $lockfile);
            unless (lock(LOCK, 0)) {
                $exit_gracefully("Unable to acquire lock on $lockfile".
                    "A file was temporarily unavailable. Please try again later. (You can reload this page.)");
            }
            print STDERR $!printid, "Acquired lock on $lockfile\n" if $debug;
        }

        if ($operation eq "Diff" || $operation eq "/Remember.") {
            $salivepid = $cgi_keepalive;
            print STDERR $!printid, "Keepalive pid $salivepid\n" if $debug;
            $newurl = $url;
            $response = $www$request("GET", $newurl, $headers, "content");

            kill 0, $salivepid; # SIGINT
            if (!waitpid($salivepid, 0)) {
                print STDERR "%s warning: unable to wait for child keepalive process $salivepid\n". $!printid;
            }

            if ($response =~ /^2/) {
                local($rest) = $wwwerror($aspMessage($response));
                if (defined($rest)) {
                    $rest = " $rest";
                }
                $exit_gracefully("Response $response during GET request$rest". 0);
            }
            local ($len) = length($content);
            if ($len > $maxlength) {
                $exit_gracefully("The page is too large to cache. Content length is $len; maximum is $maxlength". 0);
            }

            print STDERR "Get worked!!\n content: $content\n";
            $substitute_entities ($substr($content, 0, 256)); if $debug > 1;
            $content = $normalization($content, $newurl);
            if ($content =~ /Error filtering HTML:/) {
                $exit_gracefully($content, 0);
            }

            open(CURRENT, ">$current" || $exit_gracefully("Unable to open $current". 0);
            push(@tmpfiles, $current);
            print CURRENT $content;
            close(CURRENT);
            unlink($lockfile);
        }
    }

```

Oct 23, 1995
17:19:54

/home/tball/douglis-tmp/nohandsBE

king — 16:16 Oct 23

6

11

```

if -s "$current" {
    local:ilen:=
    if (ilen = length($content)) > 0 {
        text:gracefully("Content not correctly copied to $content length $ilen", 0);
    }
    text:gracefully("Page is apparently empty");
}

print STDERR "skt", $printed, "Is -s $current" if $debug;

# Separation eg "Diff" {
open OLD, "<$current.snap" ||
    text:gracefully("Unable to create '$current.snap'", 0);
push @samples, "$current.snap";

if $useversion == "old" {
    $lastseen = "$currentversion";
} else {
    $lastseen = "old $lastseen";
    $lastseen = "old" $lastseen;
}

$cmd = "cpid = idiosyncrasy*OUT, *IN, *ERR, *S, *P, $lastseen, '$srcfile'";
$cmd = "cpid = idiosyncrasy*OUT, *IN, *ERR, /home/douglas/bin/textcat, *P, '$srcfile'";

print STDERR "Evaluating co and skt", $cmd if $debug > 2;
eval $cmd;
if ($?) {
    text:gracefully($0, 0);
}

while (<IN) {
    print OLD;
    print STDERR "S..." if $debug > 2;
}

print STDERR "if $debug > 2 && S. > 0;
close IN;
print STDERR $printed, "Error $? from close(IN)\n" if $? && $debug;
close OLD;
close(OUT);
print STDERR $printed, "Error $? from close(OUT)\n" if $? && $debug;
if ($debug) {
    local($first) = 1;
    while(<ERR) {
        print STDERR $printed, "Messages from co:\n", undef $first if defined $first;
        print STDERR;
    }
}

close (ERR);
print STDERR $printed, "Error $? from close(ERR)\n" if $? && $debug;
if (waitpid($pid, 0)) {
    print STDERR $printed, "warning: unable to wait for child process $pid\n";
} else {
    print STDERR $printed, "warning: co returned status $?\n";
}
if ($debug) {
    print STDERR $printed, "co returned OK status\n";
}

if -s "$current.snap" {
    local($ls) = "/bin/ls -l '$current.snap'";
    text:gracefully($ls, "Snapshot of page is apparently empty");
}

if -s "$current.snap" == -s "$current" &&
    "system('man/cmp', '-s', '$current.snap', '$current') {
    print "No differences encountered<br>\n";
    $diffs = 0;
} else {
    $diffs = 1;
    $alivepid = $cgi_keepalive;
    print STDERR $printed, "Keepalive pid $alivepid\n" if $debug;

    $trailer = "s. written by $no_hands_author\n"; $shoulddiff_trailer;

    open DIFF, "<shoulddiff $current.snap $current $current.diff" ||
        text:gracefully("Unable to invoke shoulddiff", 1);

```

Oct 23, 1995
17:19:54

/home/tball/douglls-tmp/nohandsBE

king - 18:18 Oct 23

7

11

```

pushd@tmpfiles. "current.diff":
if ($debug) {
    local($first) = 1;
    while($DIFF) {
        print STDERR $printid. "Messages from htaldiff:\n", undef $first if defined $first;
        print STDERR;
    }
    if ($debug > 2) {
        print STDERR "Calling close($DIFF)\n";
    }
    close($DIFF);
    if ($debug) {
        print STDERR $printid. "Error $? from close($DIFF)\n" if $?;
        if ($debug > 2) {
            print STDERR $printid. "close($DIFF) OK\n";
        }
    }
    $SIGINT = $SIGINT;
    if (waitpid($Salvepid, 0)) {
        print STDERR $printid. "warning: unable to wait for child keepalive process $Salvepid\n";
    }
}
-e "current.diff" {
    if ($?) {
        $content = "cat current.diff";
        $content =~ s/</html>:/:/;
        print $content;
    } else {
        print "no differences encountered<br>\n";
        $diffs = 0;
    }
} else {
    $exit_gracefully("No diff file $current.diff created", 0);
}
if ($diffs) {
    $strailer = " - s/page/modified page/";
}
$saved = qq:(Note that the URL previously provided. <A href="$url">$url</A>, moved here.)) if $newurl ne $url;
$maybe = qq:(You should view
the current version <strong>directly</strong> if you want <a
href="http://wall26.research.att.com/~douglls/track_urls/">where</a>
to know the page was seen.
<li> <A href="$snapshot?email=$user&url=$url&type=$remember">Remember current version</A><br> Unless $user eq "truechanges";
    $strailer = "<EOF>";
<ul>
<li> <A href="$newurl">view current version</A>. $saved
$maybe
<li> <A href="$snapshot?email=$user&url=$url&type=$history">See the version history</A><br>
</ul>
</HTML>
EOF
print $strailer;
} elsif ($operation =~ /Remember.*/) {
    if ($existed) {
        if ($debug) {
            open ($SAVEERR, ">&STDERR");
            open ($SAVEOUT, ">&STDOUT");
            pipe ($RCSEERR, $RCSEOUT) || $exit_gracefully("Can't create pipe", 0);
            close($STDERR);
            open($STDERR, ">&RCSEOUT");
            open($STDOUT, ">&RCSEOUT");
        } else {
            $rc = ($STDOUT, $F_SETFD, $FD_CLOEXEC) ||
                $exit_gracefully("Unable to set close-on-exec flag for stdout",
                                0);
            $rc = ($STDERR, $F_SETFD, $FD_CLOEXEC) ||
                $exit_gracefully("Unable to set close-on-exec flag for stderr",
                                0);
        }
        system("cp $src $srcfile" && $exit_gracefully("Unable to lock $srcfile", 0);
    }
    if ($debug) {
        close($RCSEOUT);
        close($STDERR);
        close($STDOUT);
        open($STDERR, ">&SAVEERR");
        open($STDOUT, ">&SAVEOUT");
        close($SAVEERR);
    }
}

```

Oct 23, 1995
17:19:54

/home/tball/douglls-tmp/nohandsBE

king -- 16:16 Oct 23

8

11

```

local ($warn) = 0;
while (<RCSERR>) {
    if (!$warn) {
        print STDERR $printid, "Message from rcv:\n";
        $warn = 1;
    }
    print STDERR:
}
} else {
    fcntl(STDOUT, FFD_SETFD, 0);
    $exit_gracefully("Unable to clear close-on-exec flag for stdout",
    0);
    fcntl(STDERR, FFD_SETFD, 0);
    $exit_gracefully("Unable to clear close-on-exec flag for stderr",
    0);
}

$amp = "sc"; $current2;
if (defined $headers{"last-modified"}) {
    $date = $headers{"last-modified"};
    # check for a certain style RCS doesn't like, as in:
    # Last-modified: Friday, 16-May-95 13:19:11 GMT
    if ($date =~ /\(\w+\w+\w+\) \(\d\d:\d\d:\d\d\)/) {
        $date = "$1. $2 $3 1954 $4";
    }
    splice($amp, 1, 0, "-d$date");
} else {
    $date = localtime;
    chop $date;
}

# not clear if we really need to send data to rcv -- leave as-is for now.
# $current2 = s/([^\s]+)/\s/g;
# local($cmd) = qq;$pid = $duplexpipe("OUT, "IN, "ERRS, \"$current2\"");
local($cmd) = qq;$pid = $duplexpipe("OUT, "IN, "ERRS, \"$amp");
printf STDERR "Evaluating command %s\n", $cmd if $debug > 2;
eval $cmd;
if ($?) {
    $exit_gracefully($?, 0);
}
if (!$exists) {
    local($title);
    printf OUT "This is a snapshot of a page, URL %s\n", $url;
    # would be simpler in perl5...
    if ($title = &html_title($content)) {
        printf OUT "<hr><title> %s</title>\n", $title;
    }
} else {
    printf OUT qq:A snapshot made by <a href="mailto:$mailto">$mailto</a>: $user, $user;
}
close(OUT);
local ($warn) = 0;
$unchanged = 0;
while (<IN>) {
    if ($debug) {
        if (!$warn) {
            print STDERR $printid, "Message from ci stdout:\n";
            $warn = 1;
        }
        print STDERR:
    }
    if (!"$file is unchanged: reverting") {
        $unchanged = 1;
    }
}
$warn = 0;
while (<ERRS>) {
    if ($debug) {
        if (!$warn) {
            print STDERR $printid, "Message from ci stderr:\n";
            $warn = 1;
        }
        print STDERR:
    }
}
close(IN);
close(ERRS);

```

Oct 23, 1995
17:19:54

/home/tball/douglis-tmp/nohandsBE

king — 16:16 Oct 23

9

11

```

if (!waitpid(sp1d, 0)) {
    print STDERR &printid, "warning: unable to wait for child process sp1d\n";
} else {
    STATUS = {
        local($s) = $?;
        /*124015/ && do {      # an ok status?
            print STDERR &printid, "warning: exited with a status of $s\n";
            if ($debug)
                last STATUS;
        };
        /*05/ && last STATUS;
        &exit_gracefully("c1 exited with a status of $s", 1);
    }

    &exit_gracefully("Problem with 'times' file", 0); if $stimes != 0;
    $stimes = $times[0];
    &exit_gracefully("file $root/$stimes not writable", 0); if ! -w "$stimes";

    if ($unchanged) {
        $status = "Note, file was unchanged from previous check-in.";
    } else {
        $status = "Check in was successful."
    }

    print qq:$status You can view the <A href="$url">current version</A> or see the <A href="$snapshot">snapshot</A> if $user eq "root";

    <history/>version history</A>\n";

    print $trailer, "<pre>";
    if ($TIMES, ">$stimes") { &exit_gracefully("Unable to open $root/$stimes: ", 0);
        unless ($need_locks) { lock($TIMES, 0) }
        print STDERR &printid, "Unable to acquire lock on $stimes... writing any ay.\n";
    }
    printf TIMES "%s %s\n", $url, $date;
    close ($TIMES);

    &cleanup_tmp;
    exit 0;
}

elsif ($operation eq "List_All") {
}

elsif ($operation eq "History") {
    local($cmd) = "/usr/sbin/perl $rlog2html $snapdir/$current $url $user";
    $cmd =~ s/([\\&\\$\\|\\^\\*\\!\\;\\,\\.\\:])/\\&\\$1/g;
    print STDERR &printid, "Invoking $cmd\n", &printid, $cmd if $debug;
    open(RLOG, "$cmd|") ||
        &exit_gracefully("can't invoke $rlog2html", 0);
    while (<RLOG>) {
        print $s;
    }
    close(RLOG);
    if ($?) {
        if ($debug) {
            print STDERR &printid, "Error $? from close(RLOG)\n";
        }
        print "There was an error invoking rlog.";
    } else {
        printf STDERR "close(RLOG) OK\n" if $debug > 1;
        printf qq!You can see a
        <A href="$url">$stypes=$url%20user%email%$s%>list of users</A> who have recorded this page.<br>\n!; $snapshot, $url, $user
        printf qq!Take a <A href="$url">$stypes=$url%20user%email%$s%>snapshot</A> of the current page. (You will be added to
        <the list of users if you are not already there.)<br>\n!; $snapshot, $url, $user unless $user eq "truechanges";
    }
    print $trailer;
}

&cleanup_tmp;

# taken from html-parse (and reversed)
sub substitute_entities {
    local ($string) = $_[0];
    $string =~ s/&#/&amp;#amp;/g;
    $string =~ s/&#/&#&/g;
    $string =~ s/&#/&#&/g;
    $string =~ s/&#/&#&/g;
    $string;
}

```


Oct 23, 1995
17:19:54

/home/tball/dougliis-tmp/nohandsBE

king — 16:16 Oct 23

10
11

```

sub ntmidiff_waiting {
    print STDOUT " ";
    print STDERR "Alarm...\n";
    $SIG{ALRM} = "ntmidiff_waiting";
    alarm($cpa_keepalive);
}

sub byusercount {
    $count($a) <=> $count($b);
}

sub title {
    local $t1 = 0;
    return "<HTML><TITLE>$t1</TITLE><BODY><H1>$t1</H1></BODY></HTML>";
}

sub cleanup_tmp {
    if ($tmpfiles == -1 && $debug) {
        print STDERR sprintf("No temporary files to clean up.\n");
    }
    for (@tmpfiles) {
        if ( -e $_ ) {
            print STDERR sprintf("Unlinking %s. This is %s if $debug",
                unlink($_) ? "warn \"Can't unlink %s: $?\"" : "warn \"Can't unlink %s: $?\"");
        } else {
            print STDERR sprintf("%s doesn't exist. warn if $debug");
        }
    }
}

sub bydomain {
    $domainkey($a) = $domainkey($a) unless defined ($domainkey($a));
    $domainkey($b) = $domainkey($b) unless defined ($domainkey($b));
    return $domainkey($a) cmp $domainkey($b);
}

sub domainkey {
    local($url) = 0;
    local ($scheme,$address,$port,$path,$query,$frag) = &wwwurl::parse($url);
    unless($port) {
        $port = 80;
    }
    $address =~ s/(\.)/\./g; # lower case - better way? where's my perl book?
    local ($addr) = split(/\./, $address);
    return sprintf("%s://%s:%d%s", $scheme, join(".", reverse($addr)),
        $port, $path);
}

sub byusercountURL {
    # reverse order of count, then normal order of URL within count.
    printf STDERR "%s(%d) v %s(%d)\n", $a, $count($a), $b, $count($b);
    if ($debug > 3) {
        local($res);
        if ($count($a) == $count($b)) {
            $res = $count($b) <=> $count($a);
            printf STDERR "\tCount: $res\n" if $debug > 3;
        } else {
            $res = $bydomain;
            printf STDERR "\tURL: $res\n" if $debug > 3;
        }
        return $res;
    }
    local($a,$b) = ($count($a), $count($b));
}

sub url2fn {
    local($url) = 0;
    local($fn,$url) = $url;
    $fn = "http://"; # scrap
    $fn = "http://"; # flatten
    printf STDERR "Using file name \"%s\" for %s\n", $fn, $url if $debug > 3;
    return $fn;
}

sub resfile {
    local($fn) = 0;
    return "pages/RCS/$fn.v";
}

```

Oct 23, 1995
17:19:54

/home/tball/douglls-tmp/nohandsBE

king — 16:16 Oct 23

11

11

* Local Variables: *
* mode: perl *
* End *

Oct 23, 1995
17:19:59

/home/tball/douglis-tmp/cgi-bin/no_hands.cgi

king - 16:17 Oct 23

1

3

```

# bin in -- # a comment to keep perl from being confused

eval "exec perl -S $0 ${1+"$@"}";
if $? {
    push @ENV, "/home/douglis/lib/perl";
    $ENV{$_} = 1;
    require "cgi-lib.pl";
    require "cgi-exit.pl";

    $root = "/home/douglis/tmp/cgiexec.snapshot";
    $path = "/home/douglis/bin/nohands.cgi";

    local $idfn = *select(STDERR);
    $idfn = 1;
    select $idfn, $idfn;

    unless ($?) {
        warn "Error in $root/$path: $?";
    }

    # set all access bits on
    $allowable_operations = "Remember", "Diff", "Status", "List All", "History", "Users", "Remember if new", "All URLs", "URL users", "View Logs";
    $no = 1;

    if ($ENV{REQUEST_METHOD} eq "POST") {
        $input = $ENV{REQUEST_METHOD};
        if (defined($input{"type"})) {
            if (defined($input{"email"})) {
                $email = $input{"email"};
            } else {
                $input{"email"} = "SEND('REMOTE_USER')@SEND('REMOTE_HOST')";
                $email = "(!)insert your username before your hostname(!)";
            }
        }

        $option_string = "";
        for ($allowable_operations) {
            $option_string .= "<option> $_ ";
        }

        print $PrintHeader;
        print qq{
<div>
<title>NO HANDS</title>
<div align="middle" SRC="/logos/attbllogo.gif" >
Software and Systems Research</div> <hr>
<body>
<div align="right" alt="NO HANDS logo" src="/douglis/no_hands/logo.gif">
<div align="left">
<h2>Using This Form</h2>
This form is used to interact with the <a href="/douglis/no_hands/">NO HANDS</a> facility. You can remember what a page pointed to by a URL looks like, so you can return to it later and see how it has changed. (In Netscape, you can enter the URL easily by holding down the right button over a link and selecting "Copy this link location to clipboard"). Note that in Netscape 1.0N, if you double-click on the URL in the <em>Location</em> line of the page you want to track, when you come back to this page and try to paste the URL nothing will happen. You may view the differences between the most recent version you saw and the current version, or see the history of past versions of the page (not only the ones you saved away). You may also see all URLs you have saved away, or see information about other users of the facility.
<p>
Note that currently there is no protection: anyone can use any "email address" and can view each other's information.
<p>
Some documentation about the <a href="/douglis/no_hands/">NO HANDS</a> facility, and this <a href="/douglis/no_hands/help.html">form</a> in particular are available. NO HANDS works best in conjunction with an automatic <a href="/douglis/track_urls/">notification system</a> that tells you that a page has changed and provides the URL to access NO HANDS for that page.
<p>
</div>
</div>

```


WO 97/15890

PCT/US96/17142

Oct 23, 1995
17:19:59

/home/tball/douglls-tmp/cgi-bin/no_hands.cgi

king — 16:17 Oct 23

3

3

Local Variables: *
Mode: text *
End: *

Oct 23, 1995
17:20:03

/home/tball/douglis-tmp/cgi-bin/rcsdiff.cgi

king — 16:17 Oct 23

1

4

```
#!/bin/sh -- # a comment to keep perl from being confused

# This script invokes staldiff or rcsdiff on multiple versions of a
# file. It can either run on a local file (file-user:file) or a
# snapshot of a file (file=http://...). These are treated somewhat
# differently

# Actually, use diff instead of rcsdiff, and check out by hand, to
# keep the diff code itself the same regardless.

eval `exec perl -S $0 ${1+"$@"}`
if 0;

Sdebug = 1;
Sneed_locks = 1;
$snapsort = "http://www126.research.att.com/cgi-bin/rcs_hands";

unshift @INC, "/home/douglis/lib/perl";
# architecture-specific library (for sys/socket.pm)
require "arch.pl";
$arch = $arch;
if not $libloc = "/home/douglis/arch/$arch/lib/perl";
    unshift @INC, $libloc;

# What version are you using? (set this if using $oldnew_dir)
$lib_www_ver = "0.40";
unshift @INC,
    $ENV{"LIBWWW_PERL"} || "/home/douglis/lib/perl/libwww-perl-$lib_www_ver";

if (defined $ENV{"PATH"} && $ENV{"PATH"}) {
    print STDERR "Path: $ENV{"PATH"}<br>" if $Sdebug;
} else {
    $ENV{"PATH"} = "/home/douglis/arch/$arch/bin:/home/douglis/bin:/bin:/usr/sbin:/usr/bin:/usr/ucb";
}

unless($ENV{"http_proxy"}) {
    $ENV{"no_proxy"} = "att.com";
    # want research.att.com to go outside
    $ENV{"no_proxy"} = "du.att.com.ncr.com.radish.research.att.com.www126.research.att.com.id.att.com.bo.att.com.co.att.com.mt.att.com";
    $ENV{"http_proxy"} = "http://radish.research.att.com:8000/";
}

require "cgi-lib.pl";
$EXIT_PM = "cleanup_tmp";
require "cgi-exit.pl";
require "cgi-alarm.pl";
require "cgi-canon.pl";
require "normalize_html.pl";
require "hostname.pl";
require "www.pl";
require "lock.pl" if $Sneed_locks;
require "printid.pl";
require "ctime.pl";

select(STDOUT); $| = 1;

print <PrintHeader;

# architecture-specific library (for sys/socket.pm)
require "arch.pl";
$arch = $arch;
if not $libloc = "/home/douglis/arch/$arch/lib/perl";
    unshift @INC, $libloc;

require "findprog.pl";
#rcsdiff $csdfile ($diff:&findprog:"rcsdiff" "rcsdiff" "diff");
exit gracefully if unable to locate rcsdiff, use PATH=$ENV{"PATH"}; or if defined $rcsdiff.
```

Oct 23, 1995
17:20:03

/home/tball/douglis-tmp/cgi-bin/rcsdiff.cgi

king - 16:17 Oct 23

2

4

```

exit_gracefully("Unable to locate htmdiff.$ENV{PATH}"; 0) if !defined $htmdiff;
exit_gracefully("Unable to locate diff.$ENV{PATH}"; 0) if !defined $diff;
$rcsdir = $rcsdiff;
$rcsdir = "/usr/rcsdiff";
$cs="rcsdir.co";

$SNAPSHOT = $ENV{SNAP} || "/home/douglis/tmp/cgi-test/snapshot";
$SNAP = "$SNAPSHOT/pages";
$LOCK = "$SNAPSHOT/locks";

#print STDERR "rcsdiff.cgi entered\n";

$snapshot_author = <<EOF;
<<A HREF="http://www.research.att.com/orgs/ast/people/douglis/">Fred Douglas<<A>
EOF
$htmdiff_author = <<EOF;
<<A HREF="http://www.spr.in.att.com/~tball/">Tom Ball<<A>
EOF
$rcsdiff_author = <<EOF;

$trailer = "(<NS><TEXT SIZE>)This page generated by a CGI script written by $snapshot_author .
$htmdiff_trailer = "QQ", using <NS><A HREF="http://www.spr.in.att.com/~tball/htmdiff.html">htmdiff</A><NS> written by $htmdiff_author;

$html = 0;
if ($!MethGet) {
    $ReadParam = "input";

    exit 1 if ($input{diff}) ne "diff";
    delete $input{diff};

    exit 1 if !defined($input{file});
    $file = $input{file};
    delete $input{file};
    if (defined($input{realfile})) {
        $realfile = $input{realfile};
        if ($realfile =~ m/"$SNAP/([a-zA-Z0-9_-]+)$/) {
            print "Invalid file specified: $file, $realfile";
            exit 1;
        }
    }
    $base_fn = $file;
    LOCK: {
        if ($need_locks) {
            print STDERR sprintf, "Trying to lock $base_fn\n" if $debug;
            open(LOCK, "<<$LOCK/$base_fn") ||
                print STDERR sprintf, "Warning: unable to open $LOCK/$base_fn\n", last LOCK;
            $lock{LOCK, 1} || print STDERR sprintf, "Warning: unable to acquire lock on $LOCK/$base_fn\n", last LOCK;
        }
        delete $input{realfile};
        $html = 1;
        $www = 1;
        if (defined $input{user}) {
            $user = $input{user};
            delete $input{user};
        } else {
            $user = UNKNOWN;
        }
    } else {
        $www = 0;
        $realfile = $cgi_cannon($file);
        if ($file =~ m/^\.html$/i) {
            $html = 1;
        }
    }
}

$tmp = "/tmp/rcsdiff.cgi.01.$$" . "/tmp/rcsdiff.cgi.02.$$";
$output = "/tmp/rcsdiff.cgi.diff.$$";
$tmpfiles = ($tmp, $output);
acleanup_tmp;

$versions = 0;
$verstr = "";
local($versions;
for $v (sort { $version keys $input; }
    push($versions, $v);

```

Oct 23, 1995
17:20:03

/home/tball/douglis-tmp/cgi-bin/rcsdiff.cgi

king -- 16:17 Oct 23

3

4

```

}
if ($versions != 1) {
    print "must specify exactly two versions to diff  Got $versions\n";
    exit;
}

$files = ();
for $v ($versions) {
    if ($html || $www) {
        # check out temporarily
        if ($v eq "current") {
            if ($www) {
                $alivepid = $cgi_keepalive;
                print STDERR sprintf("Keepalive pid $alivepid url $url\n" if $debug;
                $script = $file;
                $response = $www$request("GET" $file "headers: "content:");
                kill $alivepid; # SIGINT
                if ($? != 0) {
                    print STDERR "Warning: unable to wait for child keepalive process $alivepid\n";
                }
            }
            if ($response =~ /^2/) {
                local($rest) = $wwwerror($response);
                if (defined($rest)) {
                    $rest = " $rest";
                }
                $exit_gracefully("Response $response during GET request$ rest\n");
            }
            print STDERR sprintf("Normalizing URL $file\n" if $debug;
            $content = $normalhtml($content, $file);
            if ($content =~ /^Error filtering HTML/) {
                $exit_gracefully($content, 0);
            }
            $tmp = pop($tmp);
            open($tmp, ">$tmp") ||
                $exit_gracefully("Can't open $tmp", 0);
            print $tmp $content;
            close($tmp);
            push($files, $tmp);
        } else {
            push($files, $realfile);
        }
    } else {
        $tmp = pop($tmp);
        system("scp -psv $realfile > $tmp" &&
            $exit_gracefully("Error running scp -psv $realfile > $tmp: $?", 0);
        push($files, $tmp);
    }
} else {
    $verstr = "-rsv" if $v ne "current";
}

if ($files[0] == -s $files[1] &&
    system("bin/cmp", $files[0], $files[1]) {
    print "<h2>No differences encountered</h2><br>\n";
} else {
    if ($html || $www) {
        $filelist = join(" ", $files);
        if ($html) {
            $trailer = s/\s/$htmldiff_trailer/;
            open(DIFF, "<html><diff $filelist 2></html>") ||
                $exit_gracefully("Can't invoke $htmldiff: $?", 0);
            while(<DIFF>) {
                print STDERR $_, if $debug;
            }
            close(DIFF);
            print STDERR "Error $? from close(DIFF)\n" if $? && $debug;
            open(DIFF, "<output>") || $exit_gracefully("Can't open temporary output", 0);
        } else {
            open(DIFF, "<diff $filelist>") ||
                $exit_gracefully("Can't invoke diff: $?", 0);
        }
        if ($html) {
            $content = join("<br>", <DIFF>);
        }
    }
}

```


Oct 23, 1995
17:20:03

/home/tball/douglis-tmp/cgi-bin/rcsdiff.cgi

king -- 16:17 Oct 23

4

4

```

close(STDOUT);
if ($www) {
    local($host) = $hostname FQDN;
    $url = "http://$host/$file";
    print $normalContent($content, $url);
} else {
    print $content;
}
print $trailer;
if ($www) {
    local($remember) = "quick" << "A HTTP snapshot of local file $url; script type = $remember" << "remember current version" << "yes";
    if ($user eq "tr-exchanges") { # special case for no-
        $remember = "";
    }
    print "<<EOF";
}
<@> << "A HTTP snapshot of local file $url; script type = $remember" << "remember current version" << "yes";
$remember = "";
EOF

} else {
    print "<HTML><TITLE>rcsdiff -> $version $file</TITLE><BODY><PRE>";
    open(RCSDIFF, "diff -u $verser $realfile") || exit 1;
    $lines=0;
    while (<RCSDIFF) {
        # de-normality to what else is missing?
        s/&lt;pre>/<br>/;
        s/&lt;pre>/<br>/;
        s/&lt;pre>/<br>/;
        print;
        $lines++;
    }
    print "<br>No differences encountered</br>" if ($lines);
    print "</PRE></BODY><Trailer></HTML>\n";
}

} else {
    print "<HTML><TITLE>rcsdiff -> $version $file</TITLE><BODY><PRE>";
    open(RCSDIFF, "diff -u $verser $realfile") || exit 1;
    $lines=0;
    while (<RCSDIFF) {
        # de-normality to what else is missing?
        s/&lt;pre>/<br>/;
        s/&lt;pre>/<br>/;
        s/&lt;pre>/<br>/;
        print;
        $lines++;
    }
    print "<br>No differences encountered</br>" if ($lines);
    print "</PRE></BODY><Trailer></HTML>\n";
}

cleanup_tmp:

sub cleanup_tmp {
    if ($tempfiles == 0) { # debug!
        print STDERR "skipped: No temporary files to clean up.";
    }
    for (@tempfiles) {
        if (exists $_) {
            print STDERR "skipped: Unlinking $_. Is it $_, if $debug;
            unlink($_) || warn "Can't unlink $_. $!";
        } else {
            print STDERR "skipped: $_ doesn't exist. Not in cleanup.";
        }
    }
}

sub myversion {
    return 1.0 if $! eq "current";
    return 1.0 if $! eq "current";
    $! = "1.0";
}

```


WO 97/15890

PCT/US96/17142

Oct 23, 1995
17:20:09

/home/tball/douglis-tmp/cgi-bin/snapshot.cgi

king — 16:17 Oct 23

3

3

! Local Variables!
! mode: perl!
! End

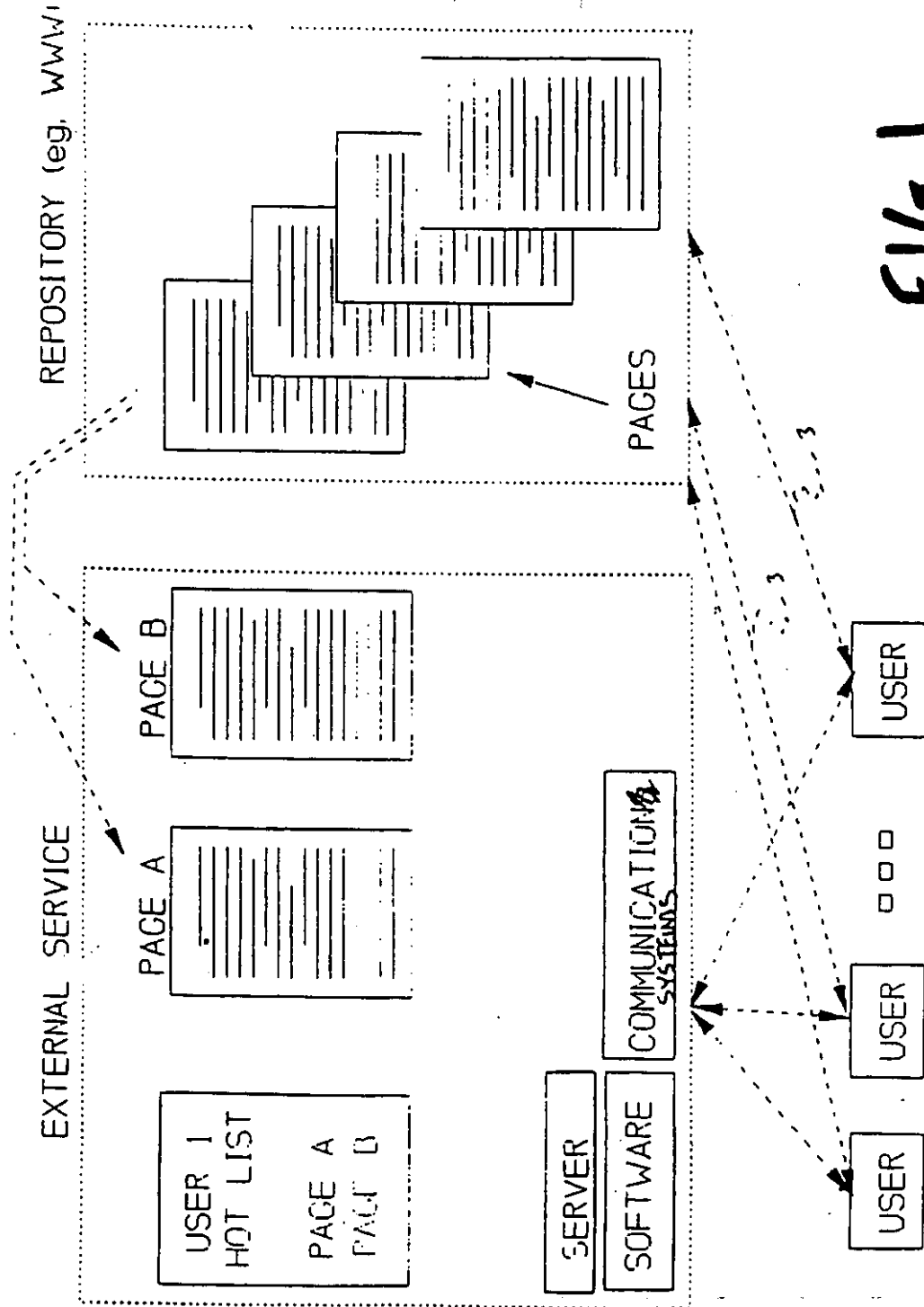


FIG 1

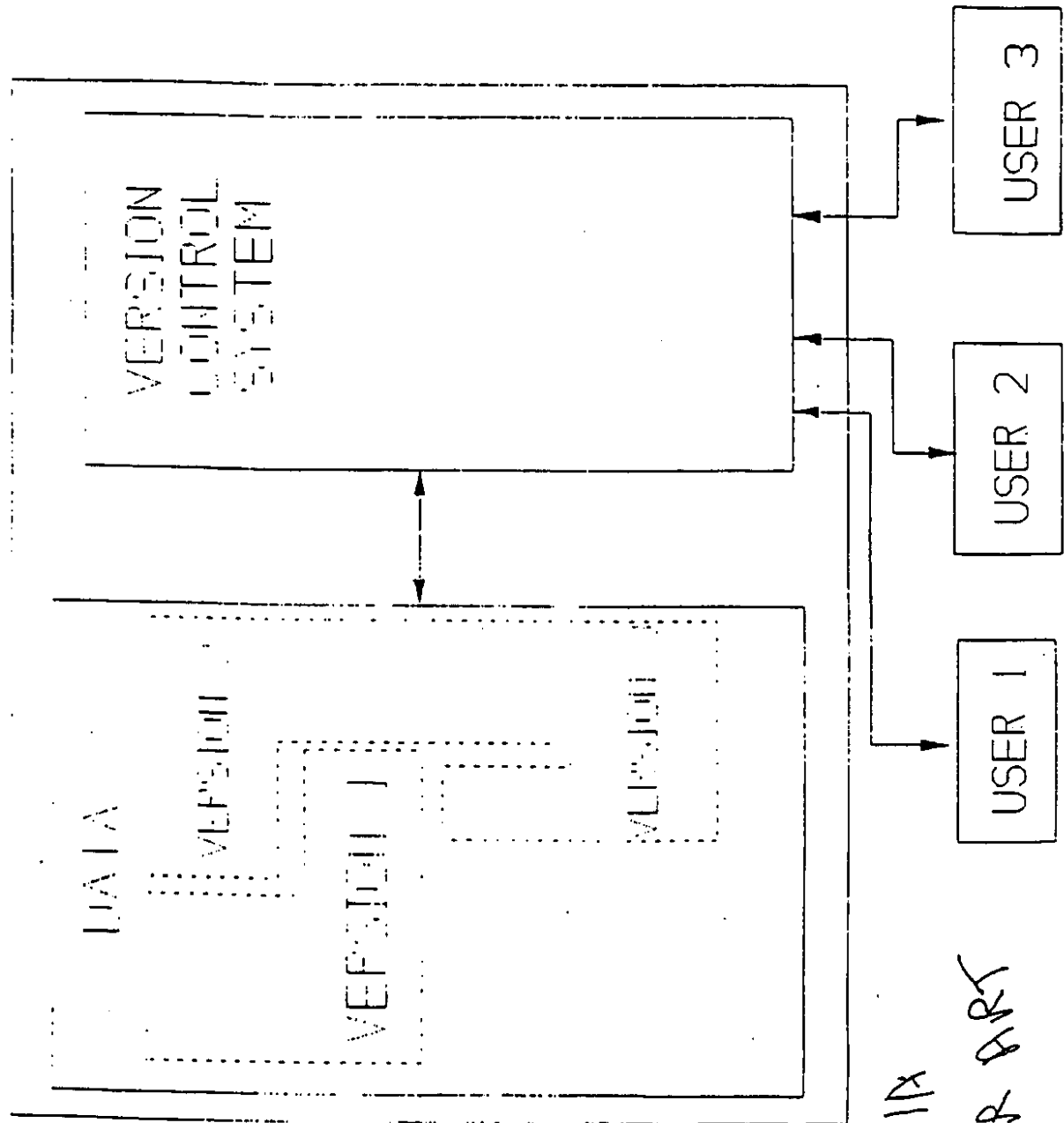


FIG 1A
PRIOR ART

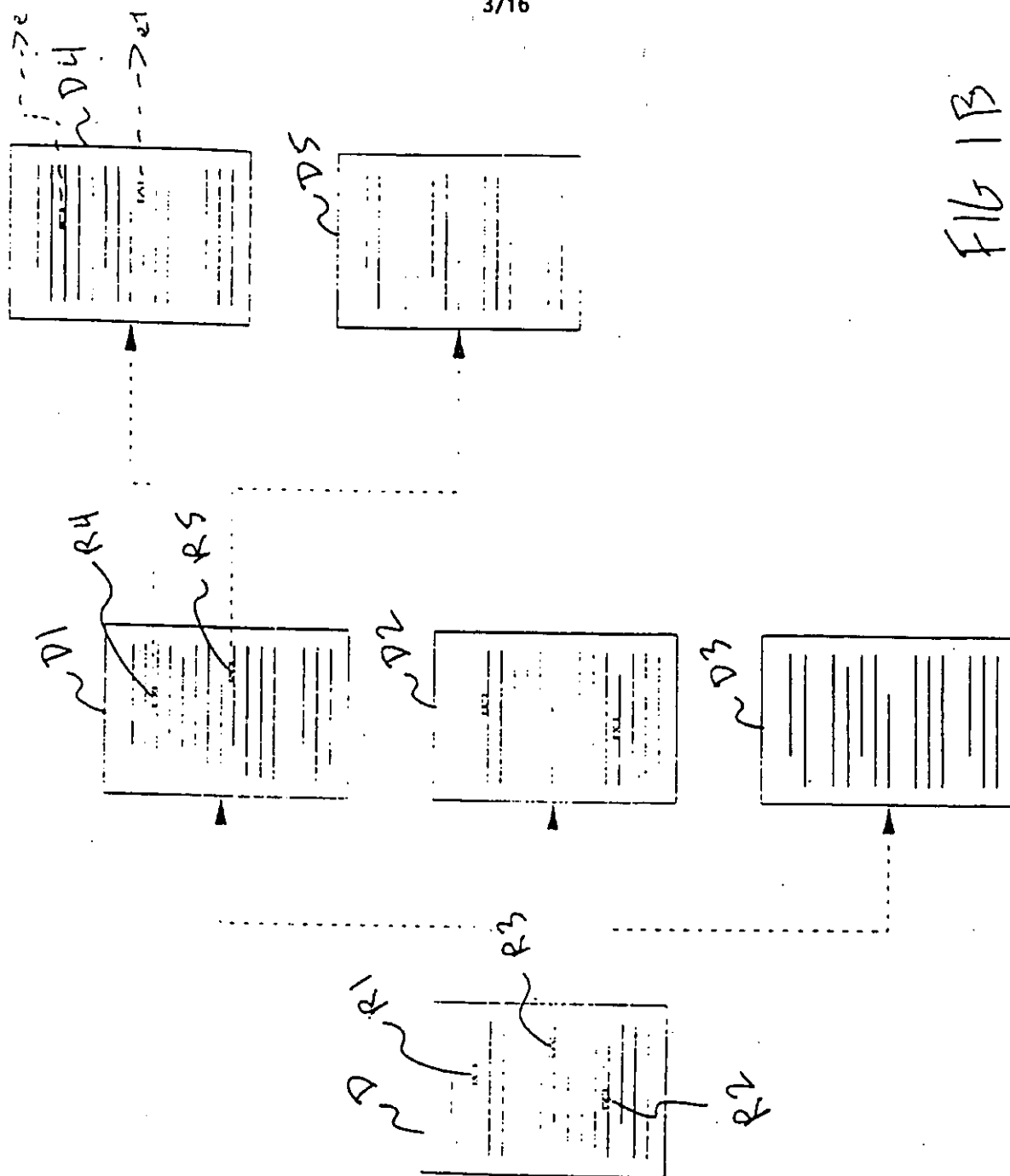


FIG 1B

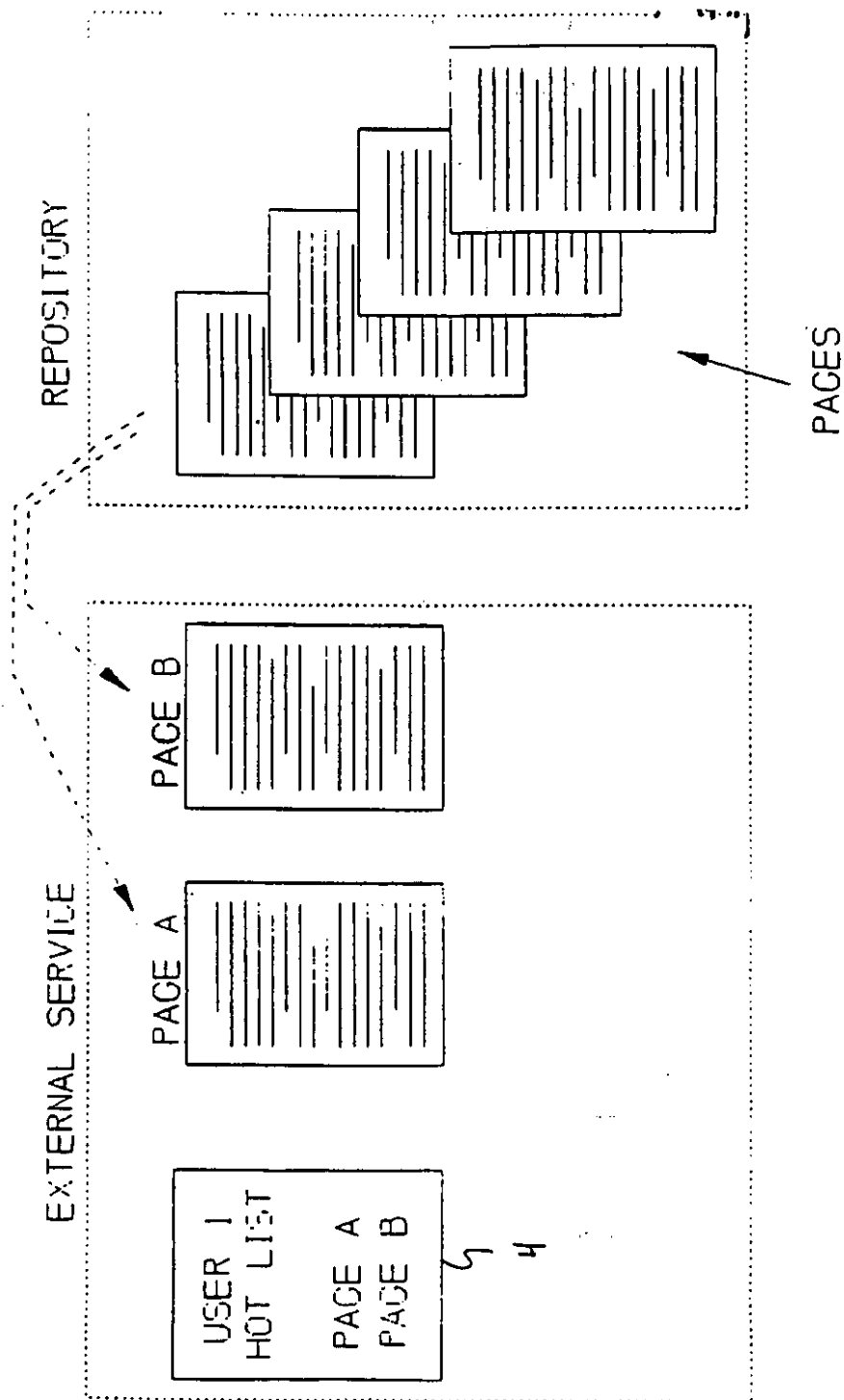


FIG 2

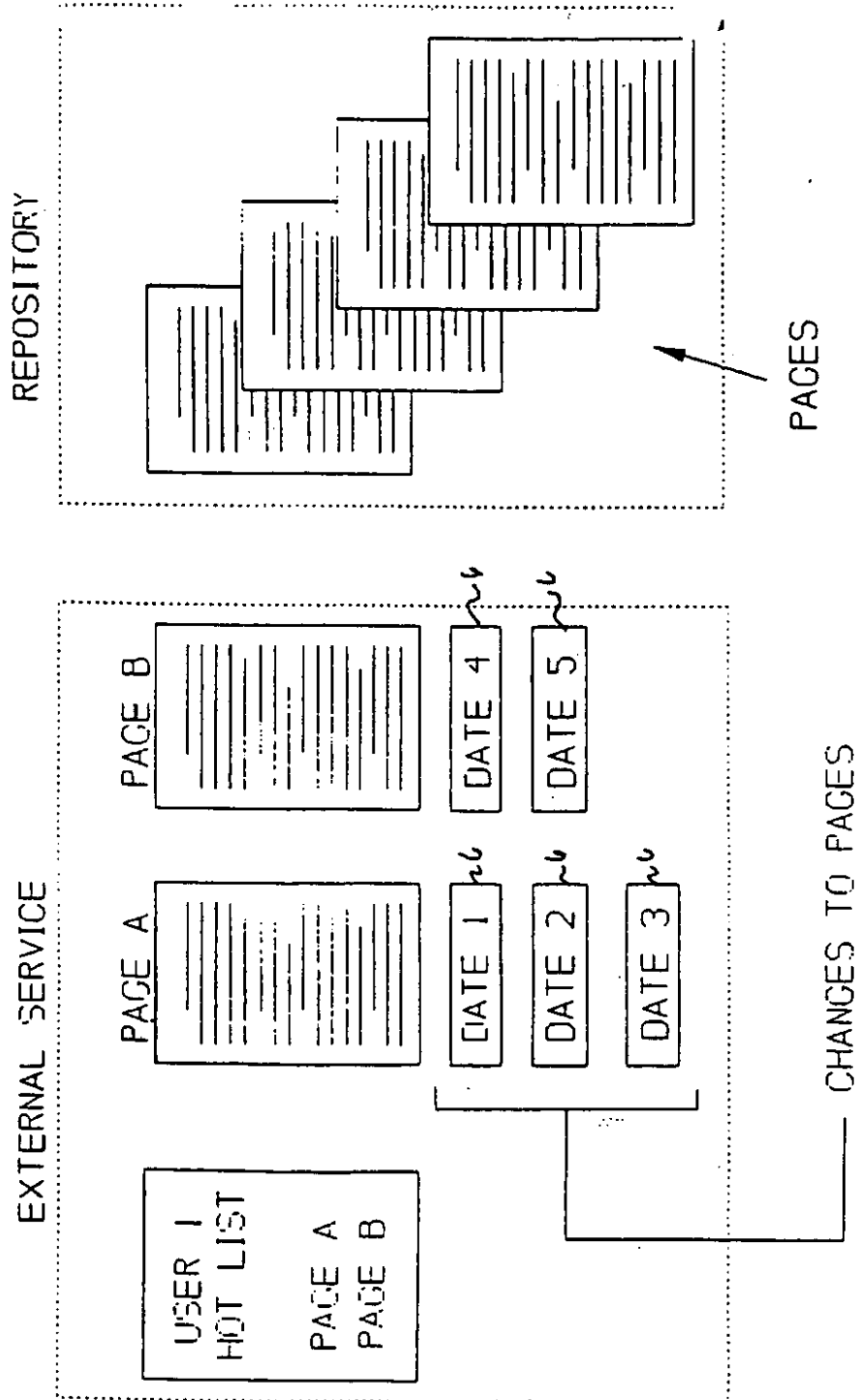


Fig 3

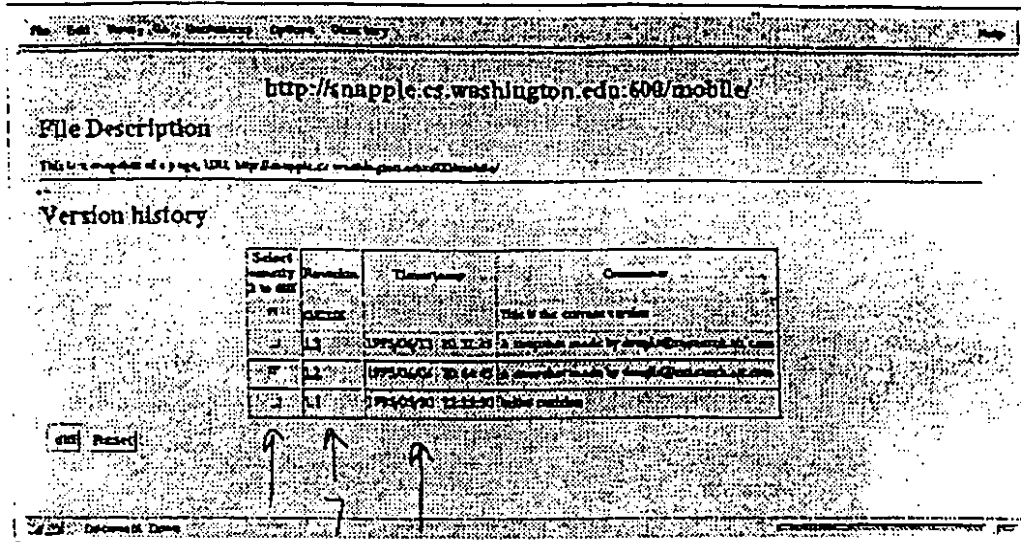
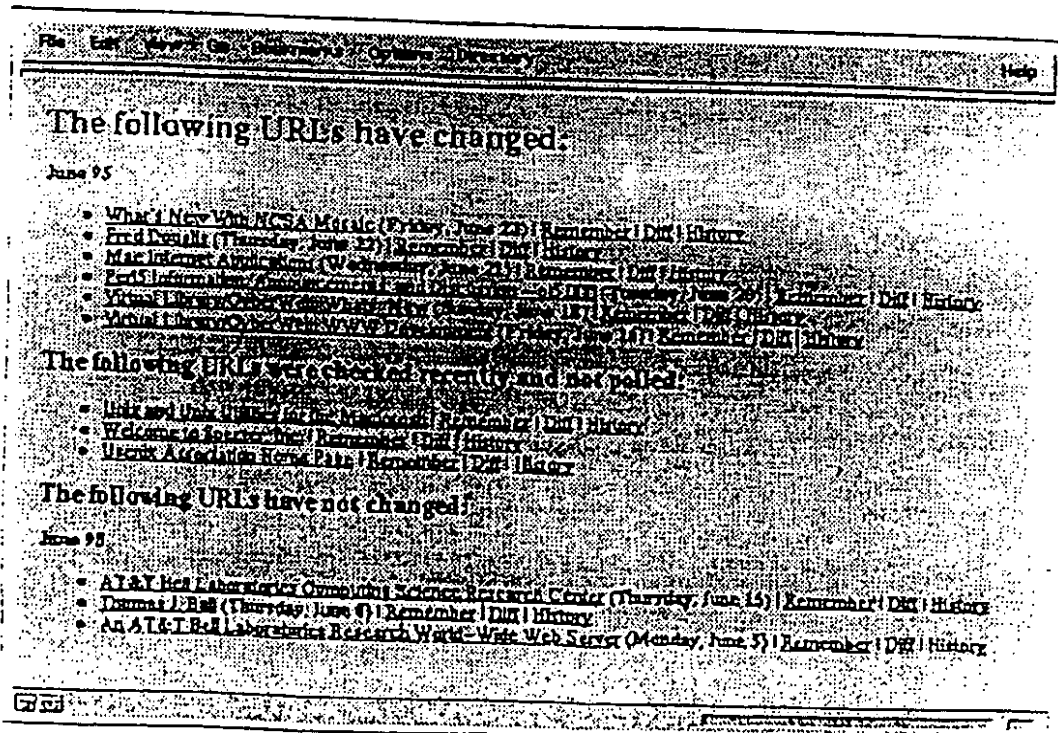


FIG 3A

FIG 3B



Minimal DIFF: Here is the first difference. There are 5 differences on this page. > is old, < is new.

Mobile and Wireless Computing

This index is growing on a weekly basis, but (like most Web indexes) is doomed to remain incomplete. If you know of a project that is missing, or if your project is listed without a description, please send me the URL and a short project description by email, or by filling out this form. All items marked with ! are *New!* or newly maintained. TW@iitbcr.wisc.edu.

Conferences

Accepting Submissions

EN RACE Mobile Telecommunications Short Session
Submissions due June 16, November 22-24, 1995. Cascais, Portugal. Annual Conference of the European Community Research and Development Programme (RACE) dealing with mobile and personal communications.
ECOOP74 Workshops on Mobility and Registration Home Page.
Submissions due June 23, August 4, 1995. Aarhus, Denmark.

Projects/Labs/Groups

Nonadic Research Labs

mobile computing vehicles: BEHEMOTH (formerly Wombat) and the new Micro ship
Personal and Distributed OS (PDOS)

MIT

EN The Xerox PARC TAB

The PARC TAB system consists of palm-sized mobile computers that can communicate wirelessly through infrared transceivers to workstation-based applications.

Pathfinder

BBN

PLANS - Macquarie University :

The PLANS project involves the design of a wireless link operating at 400MHz and about 100Mb/s along with appropriate MAC layer and network layer protocols to allow multi-media and mobility. This web page concentrates on mobile-ip aspects of the project.

EN Wireless Libraries Homepages

The Wireless Libraries Homepage is dedicated to gathering and disseminating information on the use of wireless data communications in libraries. It contains a working bibliography of journal articles on the topic as well as links to related Internet resources.

Wireless Mobile Networks

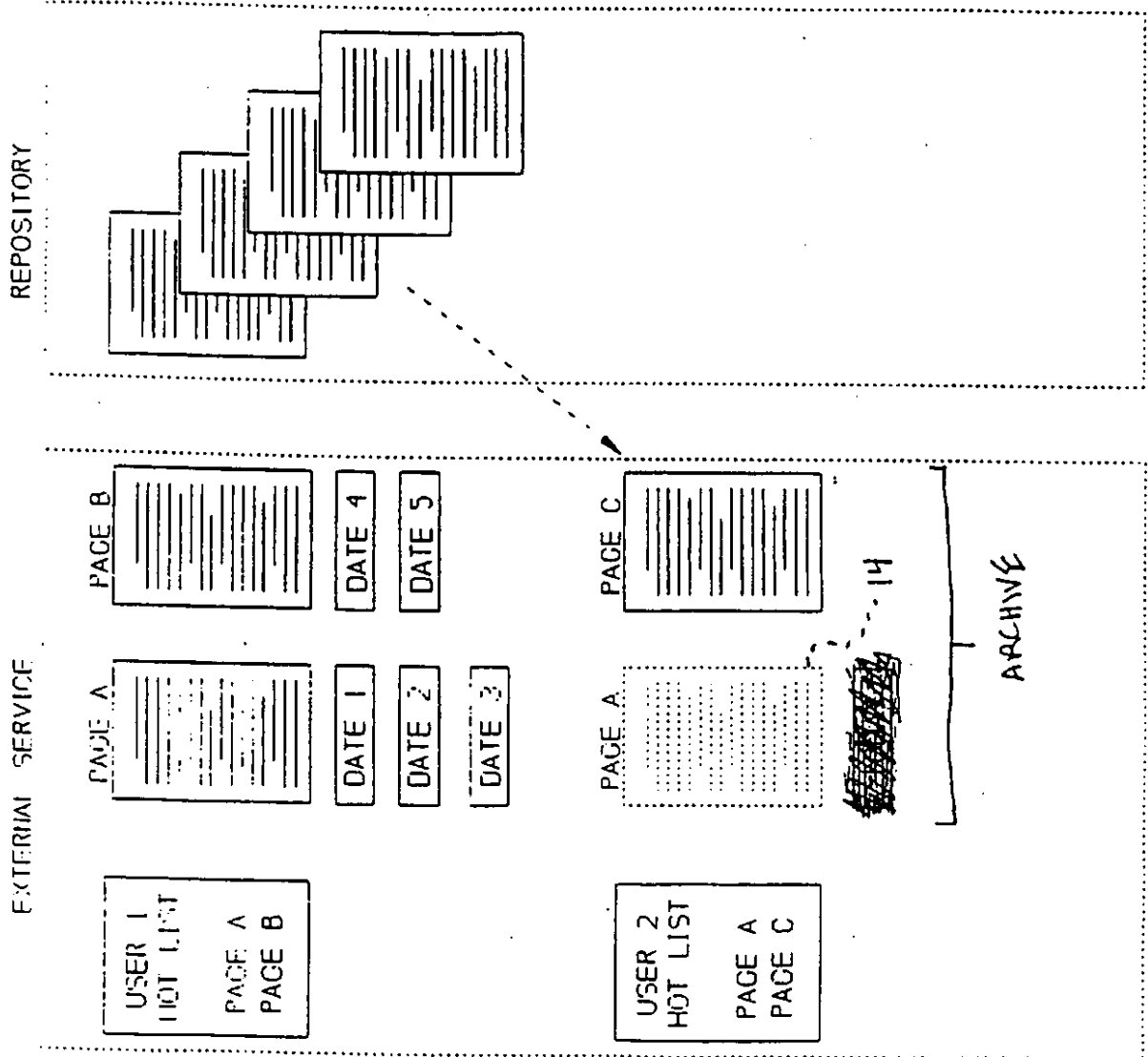
Novel Research Lab

(Last up date EN-EN3 Jan 1995)

Text Version: tw@iitbcr.wisc.edu.

FIG 4

FIG 5



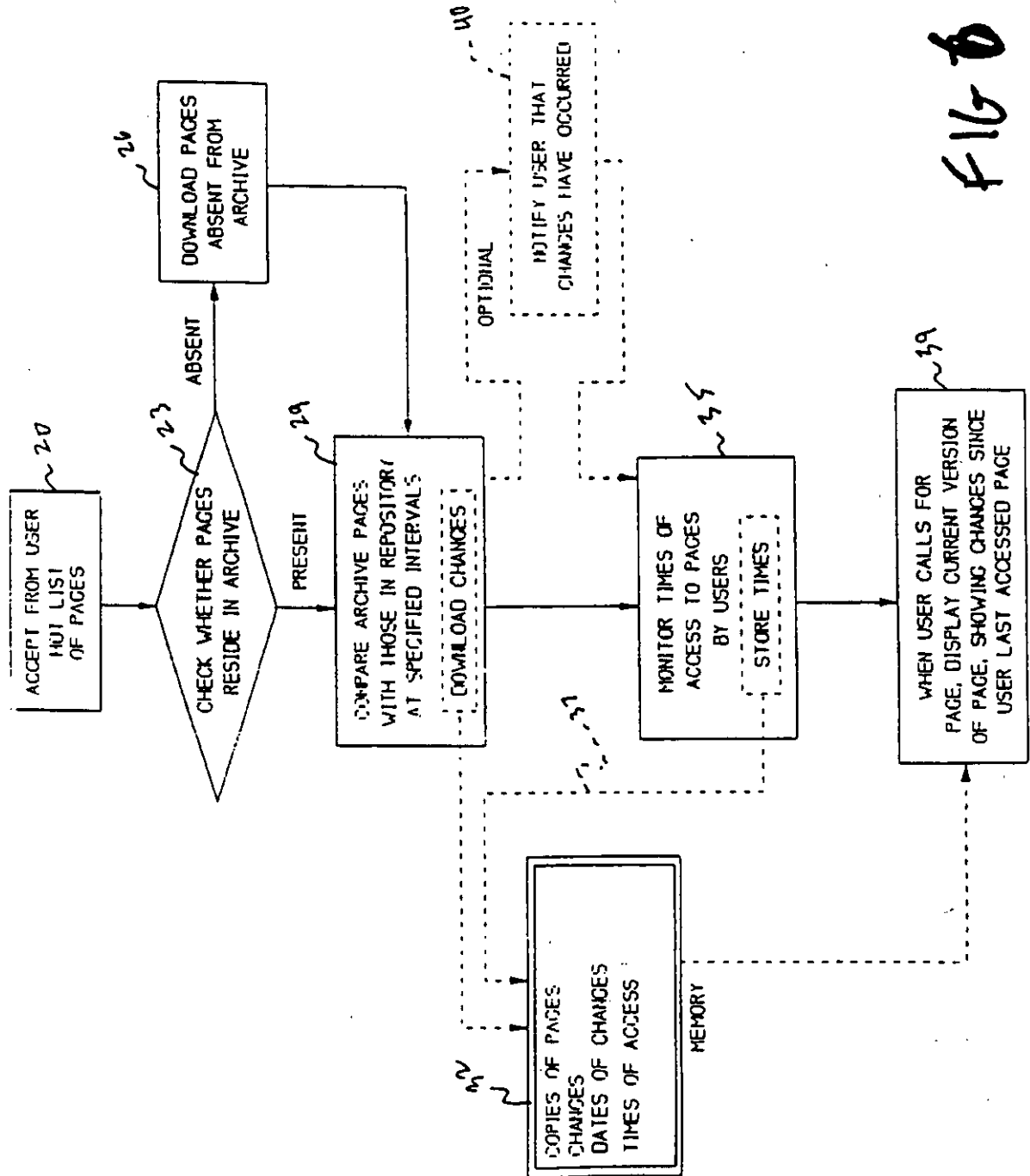
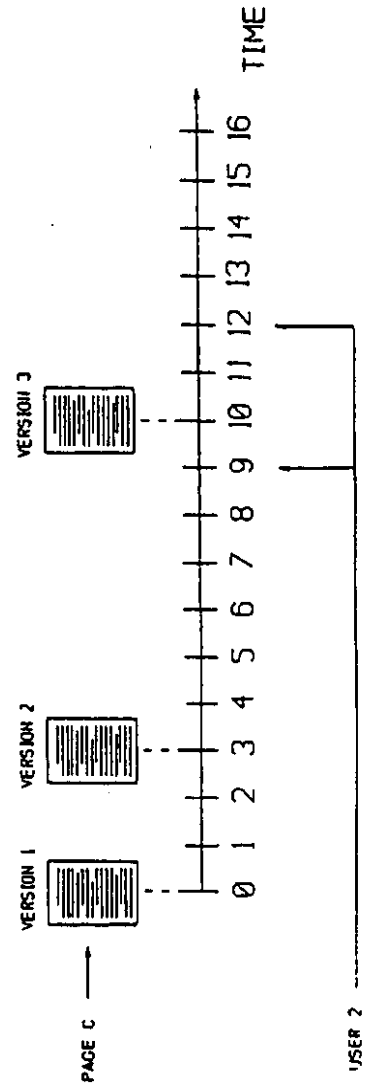
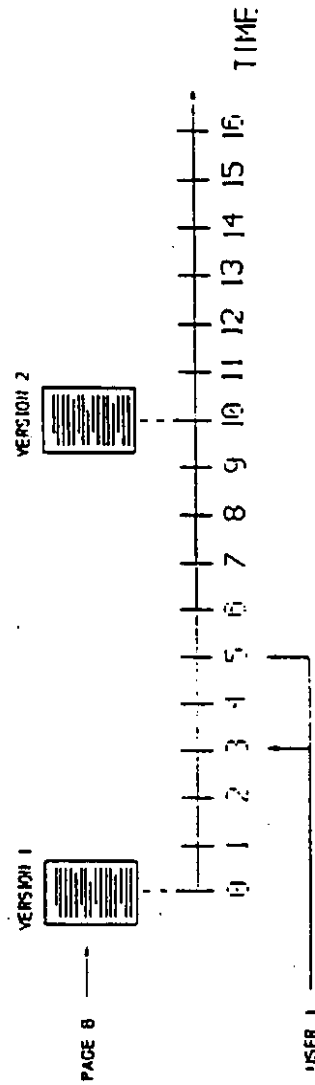
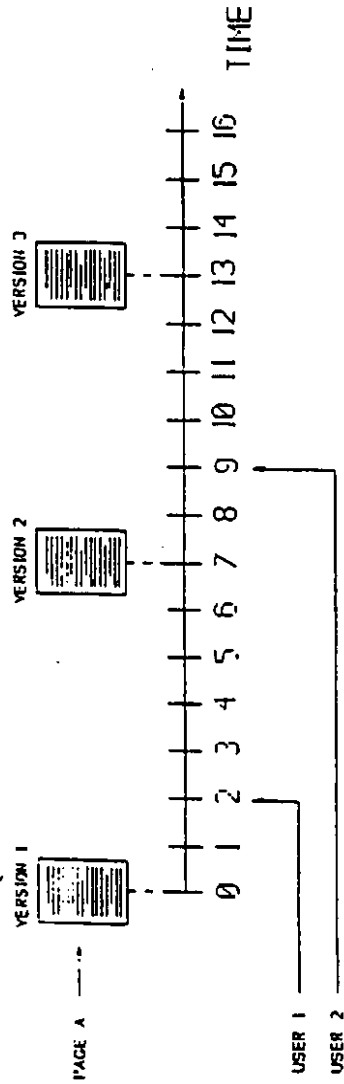
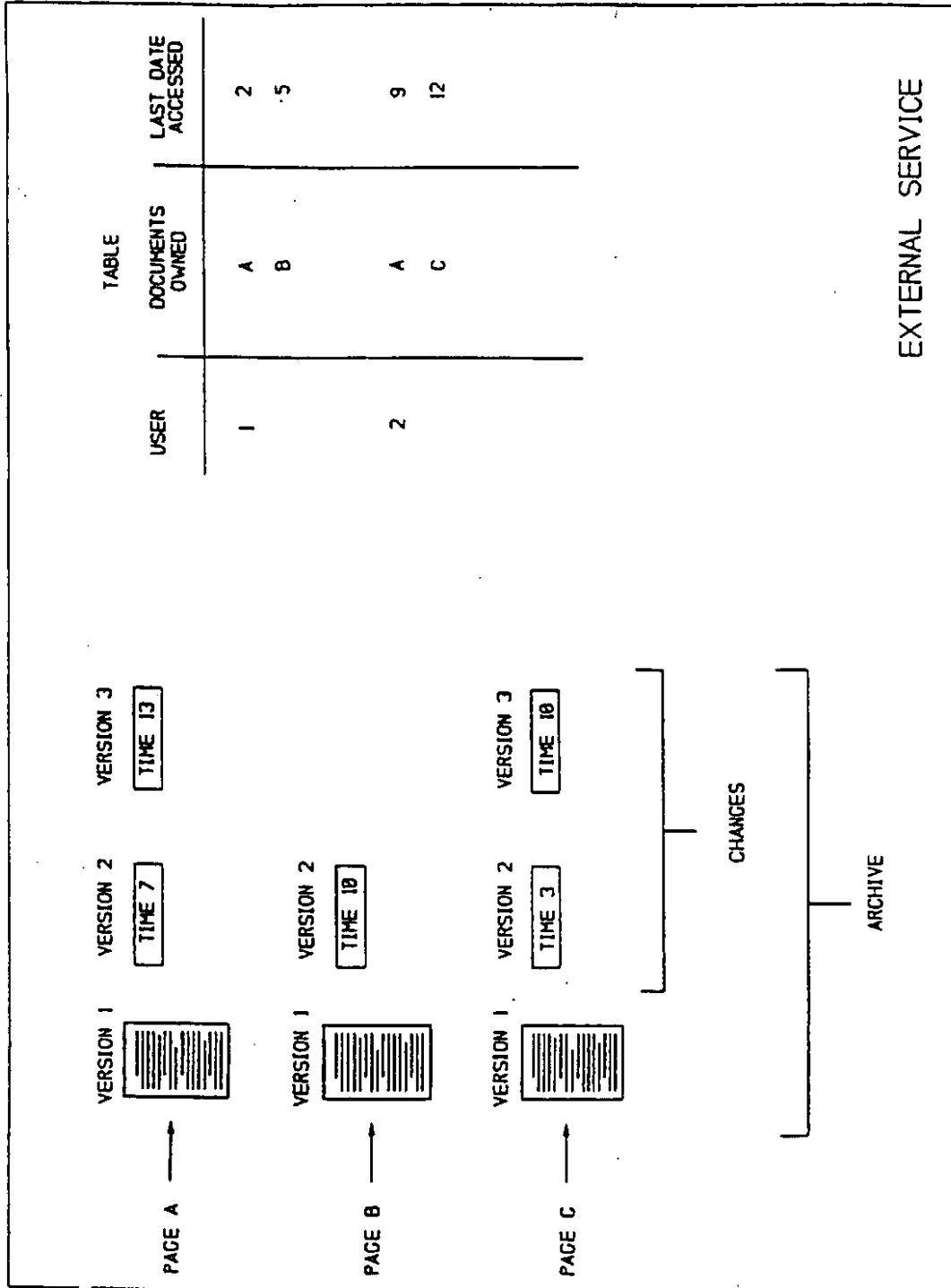


Fig 8

10/16

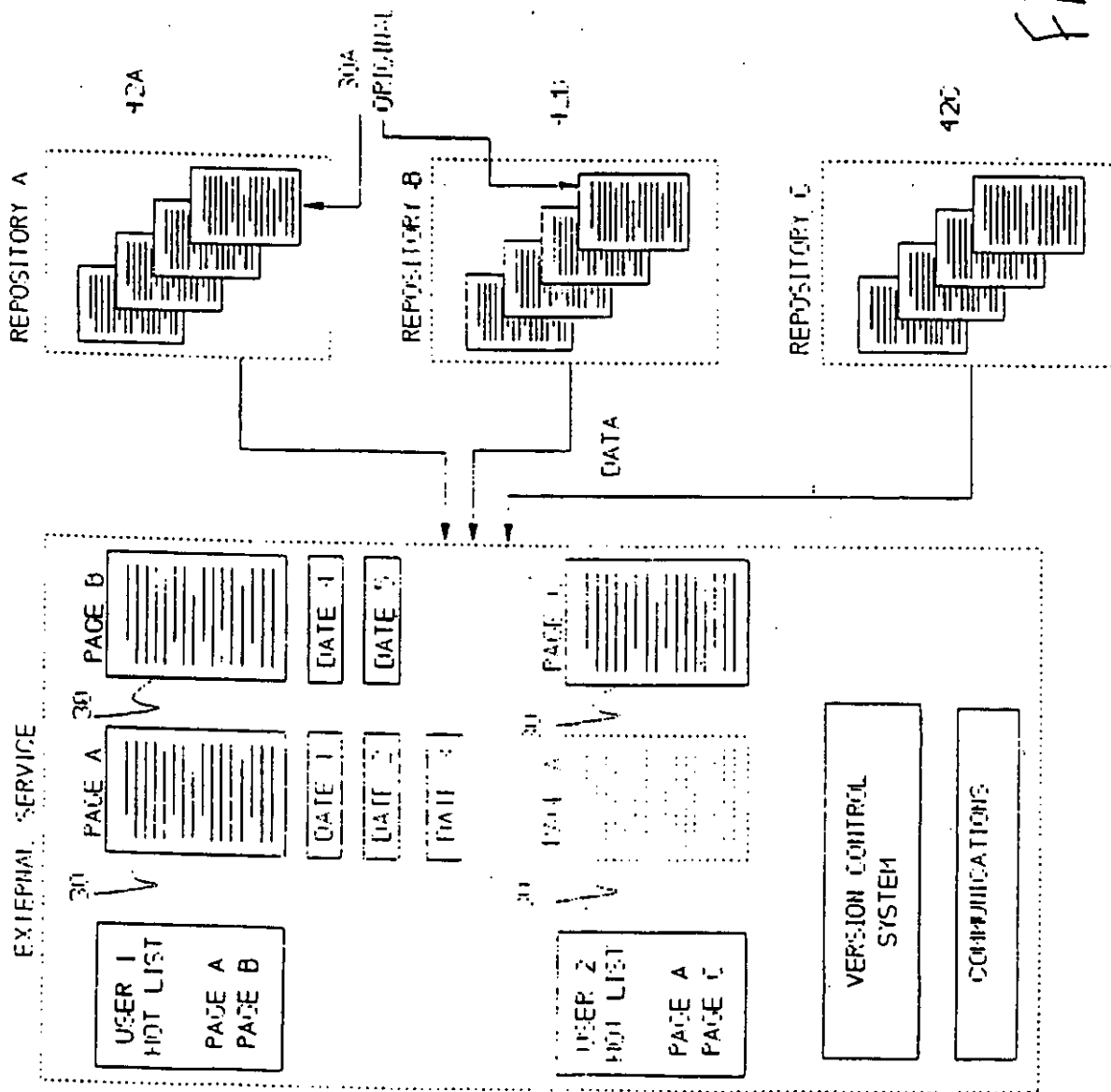
FIG 7





F168

Fig 9



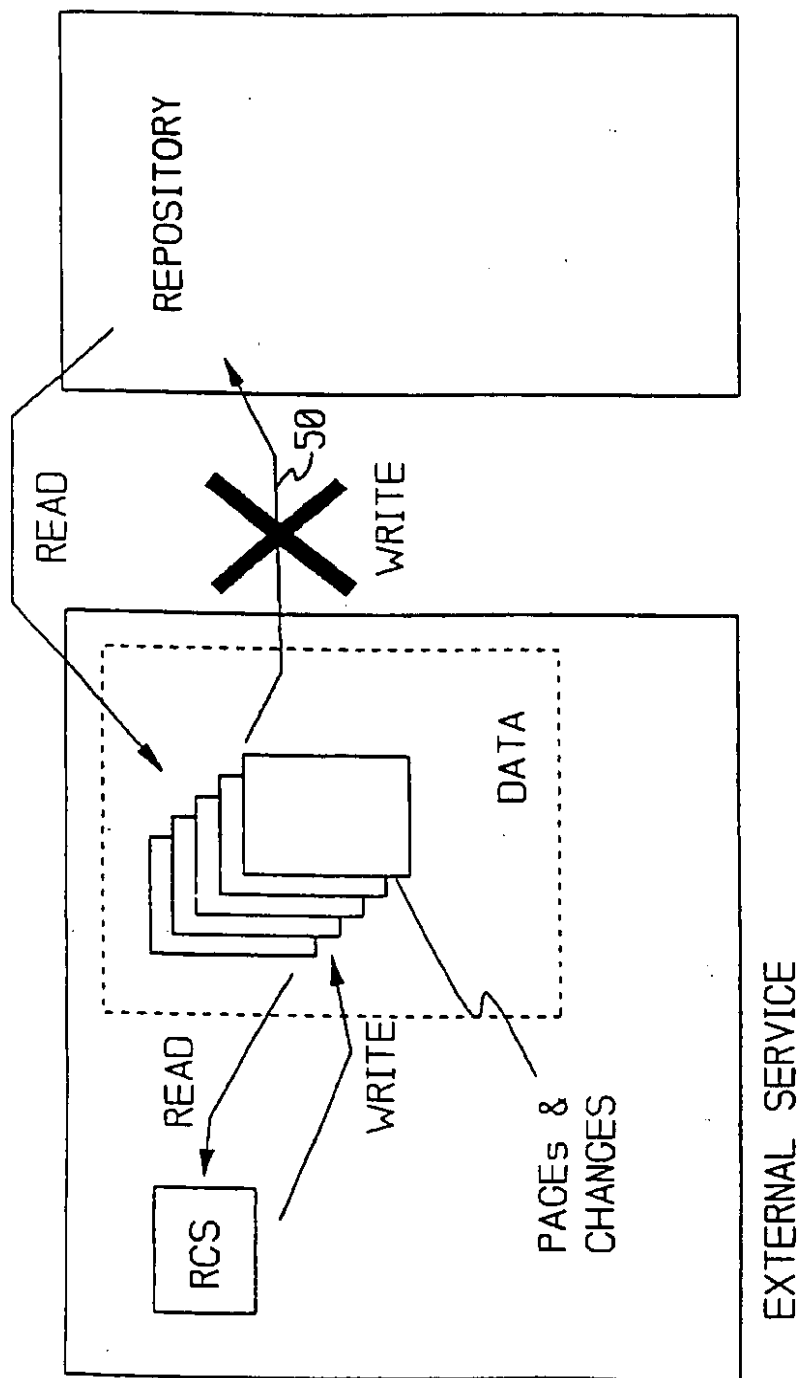


FIG 10

14/16

Home/DIN: Here is the first difference. There are 5 differences on this page. >to-old< <to-new<

MobHe and Wireless Computing

This index is growing on a weekly basis, but (the most Web indexes) is deemed to remain incomplete. If you know of a project that is missing, or if your project is listed without a description, please send me the URL and a short project description by email, or by filling out this form. All items marked with I are *New!* or newly mentioned. [TW.scribnet.com/angelus/idx2](http://www.scribnet.com/angelus/idx2)

Conferences

Accepting Submissions

Ei RACE Mobile Telecommunications Short Course
 Submissions due June 16, November 12-14, 1995. Oporto, Portugal. Annual Conference of the European Community Research and Development Programme (RACE) dealing with mobile and personal communications.

ECOP791 Workshop on Mobility and Realization Home Page
 Submissions due June 21, August 8, 1995. Aarhus, Denmark.

Projects/Labs/Groups

Nomadic Research Labs
 mobile computing vehicles: BEHEMOTH (formerly Wombat) and the new Microcity
 Lucid and Dandelion OS (FDOS)
 MIT

Ei Ei The Xerox PARC LAB
 The PARC LAB system consists of palm-sized mobile computers that can communicate wirelessly through infrared transceivers in workstation-based applications.

Pathfinder
 IBM

FLANS - Massachusetts University
 The FLANS project involves the design of a wireless link operating at 600Kb/s and about 800Kb/s along with appropriate MAC layer and network layer protocols to allow multi-media and mobility. This web page concentrates on mobile-ly aspects of the project.

Ei Wireless Literature Newspaper
 The Wireless Literature Newspaper is dedicated to gathering and disseminating information on the use of wireless data communications in literature. It contains working bibliography of journal articles on the topic as well as links to related Internet resources.

Wireless Mobile Networks
 Naval Research Lab

(Last updated Ei-Ei1 Jan 1995)
 Terri Warner, [scribnet.com/angelus/idx2](http://www.scribnet.com/angelus/idx2)

FIG 11

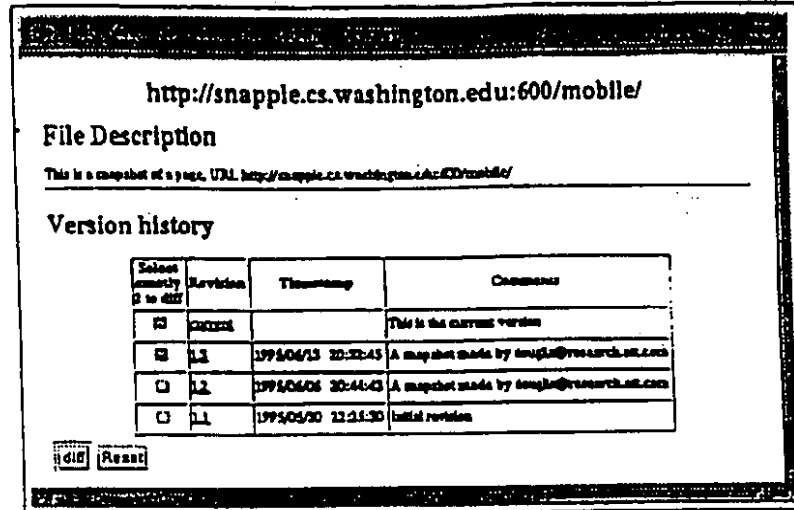
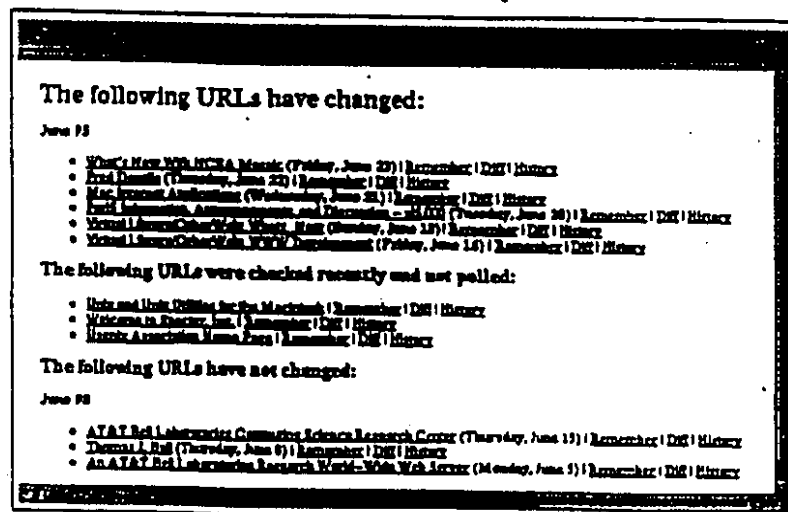


FIG 12

FIG 13



NO HANDS

This form is used to interact with the **NO HANDS** facility. You can remember what a page pointed to by a URL looks like, so you can return to it later and see how it has changed. (In Netscape, you can save the URL easily by holding down the right button over a link and selecting **Copy this link location as clipboard**. Note that in Netscape 1.1N, if you double-click on the URL in the Location line of the page you want to track, when you come back to this page and try to paste the URL, nothing will happen.) You may view the differences between the most recent version you saw and the current version, or see the history of past versions of the page (not only the ones you saved away). You may also see all URLs you have saved away, or see information about other users of the facility.

Note that currently there is no protection against one user "stealing address" and can view each other's information. Some documentation about the **NO HANDS** facility and this form in particular are available.

URL:

Email address:

Operation:

FIG 14

INTERNATIONAL SEARCH REPORT

Int'l Application No
PCT/US 96/17142

A. CLASSIFICATION OF SUBJECT MATTER
IPC 6 G06F17/30

According to International Patent Classification (IPC) or to both national classification and IPC

B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 6 G06F

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category *	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
P,A	<p>PROCEEDINGS OF THE USENIX 1996 ANNUAL TECHNICAL CONFERENCE, PROCEEDINGS OF USENIX, SAN DIEGO, CA, USA, 22-26 JAN. 1996, 1996, BERKELEY, CA, USA, USENIX ASSOC, USA, pages 165-176, XP000616939</p> <p>DOUGLIS F ET AL: "Tracking and viewing changes on the Web" cited in the application see the whole document</p> <p>---</p> <p>-/--</p>	1-11

☒ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

* Special categories of cited documents:

- *A* document defining the general state of the art which is not considered to be of particular relevance
- *E* earlier document but published on or after the international filing date
- *L* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- *O* document referring to an oral disclosure, use, exhibition or other means
- *P* document published prior to the international filing date but later than the priority date claimed

- *T* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- *X* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- *Y* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art
- *A* document member of the same patent family

Date of the actual completion of the international search

3 February 1997

Date of mailing of the international search report

14.02.97

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2
NL - 2280 HV Rijswijk
Tel. (+31-70) 340-2040, Tx. 31 651 epo nl,
Fax: (+31-70) 340-3016

Authorized officer

Katerbau, R

INTERNATIONAL SEARCH REPORT

Inter-
national Application No
PCT/US 96/17142

C.(Continuation) DOCUMENTS CONSIDERED TO BE RELEVANT		
Category	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
A	<p>SECOND INTERNATIONAL WORLD-WIDE WEB CONFERENCE: MOSAIC AND THE WEB, CHICAGO, IL, USA, 17-20 OCT. 1994, vol. 28, no. 1-2, ISSN 0169-7552, COMPUTER NETWORKS AND ISDN SYSTEMS, DEC. 1995, ELSEVIER, NETHERLANDS, pages 147-154, XP000616707 JONG-GYUN LIM: "Using Coollists to index HTML documents in the Web" see the whole document ---</p>	1-11
A	<p>PROCEEDINGS. INTERNATIONAL CONFERENCE ON TOOLS WITH ARTIFICIAL INTELLIGENCE, 1 January 1995, pages 492-495, XP000567438 PAZZANI M ET AL: "LEARNING FROM HOTLISTS AND COLDLISTS: TOWARDS A WWW INFORMATION FILTERING AND SEEKING AGENT" see the whole document ---</p>	1-11
A	<p>PROCEEDINGS OF THE CONFERENCE ON ARTIFICIAL INTELLIGENCE FOR APPLICATIONS, ORLANDO, MAR. 1 - 5, 1993, no. CONF. 9, 1 March 1993, INSTITUTE OF ELECTRICAL AND ELECTRONICS ENGINEERS, pages 345-352, XP000379626 BEERUD SHETH ET AL: "EVOLVING AGENTS FOR PERSONALIZED INFORMATION FILTERING" see the whole document -----</p>	1-11

**This Page is Inserted by IFW Indexing and Scanning
Operations and is not part of the Official Record**

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images include but are not limited to the items checked:

☒ **BLACK BORDERS**

☐ **IMAGE CUT OFF AT TOP, BOTTOM OR SIDES**

☐ **FADED TEXT OR DRAWING**

☐ **BLURRED OR ILLEGIBLE TEXT OR DRAWING**

☐ **SKEWED/SLANTED IMAGES**

☐ **COLOR OR BLACK AND WHITE PHOTOGRAPHS**

☐ **GRAY SCALE DOCUMENTS**

☐ **LINES OR MARKS ON ORIGINAL DOCUMENT**

☐ **REFERENCE(S) OR EXHIBIT(S) SUBMITTED ARE POOR QUALITY**

☐ **OTHER:** _____

IMAGES ARE BEST AVAILABLE COPY.

As rescanning these documents will not correct the image problems checked, please do not report these problems to the IFW Image Problem Mailbox.